

Particle-based likelihood inference in partially observed diffusion processes using generalised Poisson estimators

Jimmy Olsson and Jonas Ströjby

*Center of Mathematical Sciences
Lund University
Lund, Sweden*

e-mail: jimmy@maths.lth.se; strojby@maths.lth.se

Abstract: This paper concerns the use of the expectation-maximisation (EM) algorithm for inference in partially observed diffusion processes. In this context, a well known problem is that all except a few diffusion processes lack closed-form expressions of the transition densities. Thus, in order to estimate efficiently the EM intermediate quantity we construct, using novel techniques for *unbiased* estimation of diffusion transition densities, a random weight fixed-lag auxiliary particle smoother, which avoids the well known problem of particle trajectory degeneracy in the smoothing mode. The estimator is justified theoretically and demonstrated on a simulated example.

AMS 2000 subject classifications: Primary 62M09; secondary 65C05.

Keywords and phrases: auxiliary particle filter, EM algorithm, exact algorithm, generalised Poisson estimator, partially observed diffusion process, sequential Monte Carlo.

Contents

1	Introduction	2
2	Preliminaries	4
2.1	Generalised Poisson estimators	6
2.2	GPE-based particle smoothing	7
2.2.1	Convergence of the GPEPS	10
2.3	Fixed-lag smoothing	10
2.3.1	Convergence of the intermediate quantity	12
2.4	Forward-filtering backward-smoothing	13
3	Simulation study	14
3.1	Log-growth model	15
3.2	Genetics diffusion model	17
4	Conclusion	18
A	Proofs	19
A.1	Proof of Proposition 2.1	21
A.2	Proof of Theorem 2.1	26
B	More on the GPE	28
	References	30

1. Introduction

In this paper we discuss the use of *sequential Monte Carlo* (SMC) methods (alternatively termed *particle methods*) for likelihood-based inference in *partially observed diffusions* (PODs). The proposed method relies on a novel approach for estimating transition densities of diffusion processes via so-called *generalised poisson estimators* (GPEs). For the models under consideration, the likelihood function of the observed data cannot be expressed on closed-form; however, since partially observed diffusion models are, like more general latent variable models, specified using conditional dependence relations, this inference problem can be efficiently cast into the framework of the *expectation-maximisation* (EM) *algorithm* proposed by Dempster et al. (1977). When applying the EM algorithm in the POD context there are two main difficulties: firstly, in all except a few cases, the transition density of the diffusion process, and thus the complete data log-likelihood function, lacks an analytic expression; secondly, computing the *intermediate quantity* of the expectation-step involves taking expectations under the *smoothing distribution*, i.e. the conditional distribution of the hidden states at the observation time points given the observed data record, which is not—even in the case of a known transition density—available on closed-form. These two issues make, as documented by several authors, MLE-based inference in PODs very challenging. In this paper we address these problems by applying the GPE suggested (as a refinement of results obtained in Beskos et al., 2006) by Fearnhead et al. (2008) in conjunction with SMC smoothing algorithms. Unfortunately, it has been observed by several authors that using standard SMC methods in the smoothing mode may be unreliable for larger observation sample sizes n , since resampling systematically the particles leads to degeneracy of the particle paths. As a solution, we adapt the *fixed-lag smoother* proposed by Olsson et al. (2008) to the framework of PODs. This technique relies, in the spirit of Kitigawa (1998), on *forgetting properties* of the conditional hidden chain; by this is meant that the hidden chain forgets its past when evolving, backwards as well as forwards, conditionally on the given observation sequence. The constructed algorithm avoids efficiently particle trajectory degeneracy at the cost of a bias which can however be controlled by a suitable choice of the introduced lag parameter.

In order to obtain a high performance of the particle smoother it is in general necessary to propose (mutate) the particles according a kernel that takes the information provided by the current observation into account; indeed, mutating, as in the *bootstrap particle filter*, the particles “blindly” according to the dynamics of the hidden Markov chain will often lead to severe degeneracy of the particle importance weights. However, such an improved proposal strategy is not straightforwardly adopted to PODs, since computing the resulting importance weights involves computing a ratio of the transition density of the hidden diffusion process (for which a closed-form expression is missing in general) and that of the chosen proposal kernel. To cope with this, we follow Fearnhead et al. (2008) and replace each evaluation of the hidden process transition density by a draw from the GPE. Thus, the GPE serves two purposes in our algorithm

as it is used, firstly, for computing unbiased estimates of particle importance weights for a particle filter based on a proposal kernel different from the transition kernel of the hidden diffusion process and, secondly, for estimating the EM intermediate quantity itself.

The contribution of our study is fourfold, since the proposed intermediate quantity estimator

1. approximates efficiently the expectation step in a *single sweep* of the data record, yielding an algorithm with a computational complexity of order $\mathcal{O}(nN)$;
2. copes, as it is not based on any Euler discretisation or linearisation technique, efficiently with model nonlinearities;
3. has only limited computer data storage requirements, which is essential in, e.g., high frequency applications where sometimes very long measurement sequences are considered;
4. is provided with a rigorous convergence result describing its convergence to the true intermediate quantity. This result is derived via a convergence result, obtained under minimal assumptions, for the GPE-based particle smoother.

For models exhibiting poor mixing properties, in which case we cannot expect a high performance of the fixed-lag smoother, we propose an alternative algorithm where the GPE is used in conjunction with the particle-based *forward-filtering backward-smoothing* procedure proposed by Godsill et al. (2004). This scheme, which relies on a decomposition of the smoothing measure that incorporates the so-called *backward kernels* (i.e. the transition kernels of the hidden Markov chain when evolving backwards in time and conditionally on the observations) of the model, avoids particle path degeneracy completely through an additional simulation pass in the time-reversed direction. Moreover, it does not suffer from the additional, model dependent bias of the fixed-lag smoother. However, these appealing properties are obtained at the cost of a significant increase of computational work, since the complexity of the scheme in question is quadratic in the number of particles.

The paper is organised as follows: In Section 2 we recall the concept of PODs and discuss likelihood-based inference in such models via data augmentation and the EM-algorithm. GPEs are described in Section 2.1 and Section B, and Section 2.2 is devoted to SMC smoothing in general. In Sections 2.3 and 2.4 we introduce the fixed-lag smoother and the forward-filtering backward-simulation smoother, respectively; moreover, we discuss how these techniques can be adjusted to PODs using GPEs. A theoretical result describing the convergence of the fixed-lag-based estimator is found in Section 2.3.1, and in Section 3 we illustrate the method on partially observed log-growth and genetics diffusion models. In Section 4, the paper is concluded by some final conclusions and remarks. Proofs are found in Section A.

2. Preliminaries

In the following we assume that all random variables are defined on a common probability space $(\Omega, \mathcal{F}, \mathbb{P})$ and let \mathbb{E} denote expectations associated with \mathbb{P} . Denoting by $\mathbb{1}$ the indicator function and letting X be any random variable on (Ω, \mathcal{F}) , we will often make use of the short-hand notation $\mathbb{E}[X; A] = \mathbb{E}[X \mathbb{1}_A]$. Let $X \stackrel{\text{def}}{=} (X_t)_{t \geq 0}$ be continuous-time diffusion process taking values in some space $(\mathsf{X}, \mathcal{X})$, with $\mathsf{X} \subseteq \mathbb{R}^{d_X}$. More specifically, the dynamics of the process is governed by the the stochastic differential equation

$$dX_t = \mu(X_t, \theta) dt + \sigma(X_t, \theta) dW_t, \quad (2.1)$$

where $W \stackrel{\text{def}}{=} (W_t)_{t \geq 0}$ is Brownian motion. We denote by $\mathbb{W}^{(x)}$ the law of W given that $W_0 = x$ and let $(\mathcal{F}_t)_{0 \leq t}$ be the filtration generated by W . The functions $\mu(\cdot, \theta)$ and $\sigma(\cdot, \theta)$ are assumed to satisfy regularity conditions (locally Lipschitz with a linear growth bound) that guarantee a weakly unique, global solution of (2.1). We will consider a framework where the process X is only partially observed at discrete time points $(t_k)_{k \geq 0}$ through the process $Y \stackrel{\text{def}}{=} (Y_k)_{k \geq 0}$ taking values in some measurable space $(\mathsf{Y}, \mathcal{Y})$. The observations of Y are assumed to be, conditionally on the latent process X , independent and such that the conditional distribution G_θ of Y_k given X depends on X_{t_k} only. In the following we write, in order to simplify the notation, X_k instead of X_{t_k} . The dynamics of the diffusion as well as the measurement process depend on some unknown model parameter θ which is assumed to belong to some compact parameter space $\Theta \subseteq \mathbb{R}^{d_\theta}$. Our main target is to estimate θ using the maximum likelihood method. For simplicity we assume that the observation time points are *equally spaced* and denote by Q_θ and χ the transition kernel and initial distribution, respectively, of the time homogeneous Markov chain $(X_k)_{k \geq 0}$. The family $(Q_\theta(x, \cdot); x \in \mathsf{X}, \theta \in \Theta)$ is dominated by the Lebesgue-measure λ with corresponding Radon-Nikodym derivatives $(q_\theta(x, \cdot); x \in \mathsf{X}, \theta \in \Theta)$. Moreover, suppose that G_θ has a density function g_θ with respect to some measure μ on $(\mathsf{Y}, \mathcal{Y})$ such that, for $k \geq 0$,

$$\mathbb{P}(Y_k \in A | X_k) = \int_A g_\theta(X_k, y) \mu(dy), \quad A \in \mathcal{Y}.$$

Given a record $Y_{0:n} = (Y_0, Y_1, \dots, Y_n)$ (similar vector notation will be used also for other quantites) of observations, a consistent estimate of the parameter θ is ideally formed by maximising the *observed data likelihood function* $\ell_n(\theta; Y_{0:n}) \stackrel{\text{def}}{=} \log L_n(\theta; Y_{0:n})$, where

$$L_n(\theta; Y_{0:n}) \stackrel{\text{def}}{=} \int \cdots \int g_\theta(x_0, Y_0) \chi(dx_0) \prod_{k=1}^n g_\theta(x_k, Y_k) Q_\theta(x_{k-1}, dx_k),$$

A problem with this approach is that we in general cannot compute L_n on closed-form, since this involves the evaluation of a high-dimensional integral

over a complicated integrand. Since the partially observed diffusion model above is, like more general latent variable models, specified using conditional dependence relations, computation of parameter posterior distributions is facilitated significantly by maximising instead the *complete data* log-likelihood function by means of the EM algorithm: Assume that we have at hand an initial estimate θ' of the parameter vector. In the EM algorithm an improved estimate is obtained by computing and maximising the intermediate quantity $\mathcal{Q}(\theta; \cdot)$ defined by

$$\mathcal{Q}_n(\theta; \theta') \stackrel{\text{def}}{=} \mathbb{E}_{\theta'} \left[\sum_{k=0}^{n-1} \log q_{\theta}(X_k, X_{k+1}) \middle| Y_{0:n} \right] + \mathbb{E}_{\theta'} \left[\sum_{k=0}^n \log g_{\theta}(X_k, Y_k) \middle| Y_{0:n} \right]. \quad (2.2)$$

Here we have written $\mathbb{E}_{\theta'}$ to stress that the expectations are taken under the dynamics determined by the initial parameter θ' . Under weak assumptions, repeating recursively this procedure yields a sequence of parameter estimates that converges to a stationary point θ^* of the observed data log-likelihood (Wu, 1983). As clear from (2.2), computing \mathcal{Q}_n requires the computation of expected values under the smoothing distribution, i.e. the distribution of the state sequence $X_{0:n}$ conditionally on the observations $Y_{0:n}$, given by, for $A \in \mathcal{X}^{\otimes(n+1)}$,

$$\phi_n(A; \theta) \stackrel{\text{def}}{=} \frac{\int \cdots \int_A g_{\theta}(x_0, Y_0) \chi(\mathrm{d}x_0) \prod_{k=1}^n g_{\theta}(x_k, Y_k) Q_{\theta}(x_{k-1}, \mathrm{d}x_k)}{L_n(\theta; Y_{0:n})}. \quad (2.3)$$

Of special interest is the *filter distribution*, i.e. the distribution of X_n conditionally on $Y_{0:n}$, given by the restriction $\phi_{n|n}(A) \stackrel{\text{def}}{=} \phi_n(\mathbf{X}^n \times A)$, $A \in \mathcal{X}$, of the smoothing distribution to the last component. It is easily shown that the flow $(\phi_k)_{k=0}^{\infty}$ satisfies the well-known *forward smoothing recursion*

$$\phi_{k+1}(A; \theta) = \frac{L_k(\theta; Y_{0:k})}{L_{k+1}(\theta; Y_{0:k+1})} \iint_A g_{\theta}(x_{k+1}, Y_{k+1}) Q_{\theta}(x_k, \mathrm{d}x_{k+1}) \phi_k(\mathrm{d}x_{0:k}; \theta), \quad (2.4)$$

where $A \in \mathcal{X}^{\otimes(k+2)}$. By introducing the (non-Markovian) transition kernel

$$L_k(x_k, A; \theta) \stackrel{\text{def}}{=} \int_A g_{\theta}(x_{k+1}, Y_{k+1}) Q_{\theta}(x_k, \mathrm{d}x_{k+1}),$$

for $x_k \in \mathbf{X}$ and $A \in \mathcal{X}$, we may rewrite the recursion (2.4) as

$$\phi_{k+1}(A; \theta) = \frac{\iint_A L_k(x_k, \mathrm{d}x_{k+1}; \theta) \phi_k(\mathrm{d}x_{0:k}; \theta)}{\iint L_k(x_k, \mathrm{d}x_{k+1}; \theta) \phi_k(\mathrm{d}x_{0:k}; \theta)}. \quad (2.5)$$

Here the normalised (Markovian) kernel $L_k(x_k, A; \theta)/L_k(x_k, \mathbf{X}; \theta)$ is the so-called *optimal kernel* describing the distribution of X_{k+1} given $X_k = x_k$ and the new observation Y_{k+1} .

In general, a closed-form solution of the recursion (2.4) is not available. A standard approach is thus to apply some SMC smoothing algorithm (described in in Section 2.2) to approximate the expectations in (2.2). Unfortunately, both the SMC smoother itself as well as the intermediate quantity

(2.2) call for the transition density q_θ , which is usually unknown except in a few special cases. Nevertheless, results obtained by Beskos et al. (2006) and Fearnhead et al. (2008) offer a method for estimating this density *without bias*. A full treatment of this technique—which is a key ingredient of the estimation technique proposed here—is beyond the scope of this paper; nevertheless, the main framework and assumptions are described briefly in the next section. In addition, some more details can be found in Appendix B.

2.1. Generalised Poisson estimators

Define the function

$$\eta(\cdot, \theta) : u \mapsto \int^u \frac{1}{\sigma(v, \theta)} dv ,$$

and set $\tilde{X}_t \stackrel{\text{def}}{=} \eta(X_t, \theta)$. Denote by f^\leftarrow the inverse of any invertable function f . By applying Itô's formula we obtain the stochastic differential equation

$$d\tilde{X}_t = \alpha(\tilde{X}_t, \theta) dt + dW_t , \quad (2.6)$$

where

$$\alpha(u, \theta) \stackrel{\text{def}}{=} \frac{\mu\{\eta^\leftarrow(u, \theta), \theta\}}{\sigma\{\eta^\leftarrow(u, \theta), \theta\}} + \frac{1}{2} \sigma'\{\eta^\leftarrow(u, \theta), \theta\} ,$$

for the transformed process $\tilde{X} \stackrel{\text{def}}{=} (\tilde{X}_t)_{t \geq 0}$. Using again the notation $\tilde{X}_k = \tilde{X}_{t_k}$, let \tilde{q}_θ be the transition density (with respect to the Lebesgue measure λ) of $(\tilde{X}_k)_{k \geq 0}$. Then, straightforwardly,

$$q_\theta(x, x') = \tilde{q}_\theta(x, x') |\eta'(x', \theta)| . \quad (2.7)$$

Assume the following:

(A1) The process $(M_t)_{t \geq 0}$, with

$$M_t \stackrel{\text{def}}{=} \exp \left(\int_0^t \alpha(\tilde{X}_s, \theta) d\tilde{X}_s + \int_0^t \alpha^2(\tilde{X}_s, \theta) ds \right) ,$$

is a martingale with respect to $\mathbb{W}^{(x)}$;

(A2) $\alpha(\cdot, \theta)$ is continuously differentiable;

(A3) $\alpha^2(\cdot, \theta) + \alpha'(\cdot, \theta)$ is bounded from below by some function $l(\theta)$.

Under these conditions, the GPE approach developed by Fearnhead et al. (2008) makes it possible to generate random variables $\tilde{V}_\theta(x, x')$ with $\mathbb{E}\tilde{V}_\theta(x, x') = \tilde{q}_\theta(x, x')$ for any $(x, x') \in \mathbb{X}^2$, i.e. $\tilde{V}_\theta(x, x')$ estimates the transition density \tilde{q}_θ without any bias, for a large class of diffusions of type (2.6). Then, letting $V_\theta(x, x') \stackrel{\text{def}}{=} \tilde{V}_\theta(x, x') |\eta'(x', \theta)|$ yields, using (2.7), $\mathbb{E}V_\theta(x, x') = q_\theta(x, x')$. A full description of GPEs is beyond the scope of this paper; however, its main features are discussed in Appendix B. In this paper we represent the GPE by a kernel P_θ

in sense that $V_\theta(x, x') \sim P_\theta(x, x', \cdot)$. Similarly, using the related *exact algorithm* developed by Beskos et al. (2006), it is possible to construct a kernel \bar{P}_θ such that $\mathbb{E}\bar{V}_\theta(x, x', \theta) = \log q_\theta(x, x')$ for draws $\bar{V}_\theta(x, x', \theta) \sim \bar{P}_\theta(x, x', \cdot)$. Appealingly, it is in many cases (see Section 3 for examples) possible to construct P_θ and \bar{P}_θ such that the functions $\theta \mapsto V_\theta(x, x')(\omega)$ and $\theta \mapsto \bar{V}_\theta(x, x')(\omega)$ are continuous for any fixed outcome $\omega \in \Omega$, yielding unbiased estimates of q_θ and $\log q_\theta$ for all $\theta \in \Theta$ simultaneously. This useful property makes, as we will see, the GPE approach well suited to numerical (log-)likelihood function optimisation.

2.2. GPE-based particle smoothing

Since we in this part deal with the problem of sampling $\phi_k(\cdot; \theta)$ for a given *fixed* parameter value, we will throughout this section expunge θ from the notation. To begin with, we assume that we know the transition kernel density q .

In order to describe precisely how SMC methods may be used for producing approximate solutions to the smoothing recursion (2.4), we suppose that we are given a weighted sample $(\xi_{0:k|k}^i, \omega_k^i)_{i=1}^N$ of particle and associated weights, each particle $\xi_{0:k|k}^i = (\xi_{1|k}^i, \dots, \xi_{k|k}^i)$ being a random variable in \mathbf{X}^{k+1} , approximating ϕ_k in the sense that

$$\phi_k^N(f) \stackrel{\text{def}}{=} (\Omega_k^N)^{-1} \sum_{i=1}^N \omega_k^i f(\xi_{0:k|k}^i) \approx \phi_k(f), \quad (2.8)$$

where $\Omega_k^N \stackrel{\text{def}}{=} \sum_{\ell=1}^N \omega_k^\ell$, for a large class of estimand functions f on \mathbf{X}^{k+1} . Now, in order to form an updated particle sample approximating ϕ_{k+1} , as a new observation Y_{k+1} becomes available, a natural approach is to replace ϕ_k in (2.5) by its particle approximation. This yields the mixture (recall the notation δ_a for a Dirac mass located at a)

$$\bar{\phi}_{k+1}^N(A) \stackrel{\text{def}}{=} \sum_{i=1}^N \frac{\omega_k^i L_k(\xi_{k|k}^i, \mathbf{X})}{\sum_{\ell=1}^N \omega_k^\ell L_k(\xi_{k|k}^\ell, \mathbf{X})} \int_A \frac{L_k(\xi_{k|k}^i, dx_{k+1})}{L_k(\xi_{k|k}^i, \mathbf{X})} \delta_{\xi_{0:k|k}^i}(dx_{0:k}),$$

for $A \in \mathcal{X}^{\otimes(k+2)}$. Now, the aim is to simulate a new set of particles from $\bar{\phi}_{k+1}^N$ and repeat this recursively to obtain particle samples approximating the smoothing distributions at all time steps. However, since we in general cannot neither simulate draws from the optimal kernel nor compute the mixture weights $L_k(\xi_{k|k}^i, \mathbf{X})$, we apply importance sampling and draw new particles from the instrumental mixture distribution

$$\pi_{k+1}^N(A) \stackrel{\text{def}}{=} \sum_{i=1}^N \frac{\omega_k^i \psi_k^i}{\sum_{\ell=1}^N \omega_k^\ell \psi_k^\ell} \int_A \delta_{\xi_{0:k|k}^i}(dx_{0:k}) R_k(\xi_{k|k}^i, dx_{k+1}),$$

for $A \in \mathcal{X}^{\otimes(k+2)}$, where R_k is a Markovian proposal kernel and $(\psi_k^i)_{i=1}^N$ are positive numbers referred to as *adjustment multiplier weights*. We will from now

on assume that $\psi_k^i = \Psi_k(\xi_{0:k|k}^i)$ for some nonnegative function $\Psi_k : \mathbf{X}^{k+1} \rightarrow \mathbb{R}^+$ and that each kernel R_k has a density r_k with respect to λ . Simulating a particle $\xi_{0:k+1|k+1}^i$ from π_{k+1}^N is easily done by, firstly, drawing, according to the probability distribution proportional to $(\omega_k^i \psi_k^i)_{i=1}^N$, a mixture component (or ancestor) index I_k^i among $\{1, \dots, N\}$ and, secondly, extending the selected ancestor with a draw from the proposal kernel, i.e. letting $\xi_{0:k+1|k+1}^i \stackrel{\text{def}}{=} (\xi_{0:k|k}^{I_k^i}, \xi_{k+1|k+1}^i)$ with $\xi_{k+1|k+1}^i \sim R_k(\xi_{k|k}^{I_k^i}, \cdot)$. After this, the drawn particle is assigned the importance weight

$$\omega_{k+1}^i \stackrel{\text{def}}{=} \Phi_{k+1} \left(\xi_{0:k+1|k+1}^i \right), \quad (2.9)$$

where, for $x_{0:k+1} \in \mathbf{X}^{k+2}$,

$$\Phi_{k+1}(x_{0:k+1}) \stackrel{\text{def}}{=} g(x_{k+1}, Y_{k+1}) \Psi_k^{-1}(x_{0:k}) \frac{q(x_k, x_{k+1})}{r_k(x_k, x_{k+1})},$$

implying $\omega_{k+1}^i \propto d\bar{\phi}_{k+1}^N / d\pi_{k+1}^N(\xi_{0:k+1|k+1}^i)$. Finally, the weighted particle sample formed by the updated particles and weights is returned as an approximation of ϕ_{k+1} . Moreover, since the filter distribution is the marginal of the smoothing distribution with respect to the last component, an estimate of $\phi_{k+1|k+1}$ is formed by the marginal sample $(\xi_{k+1|k+1}^i, \omega_{k+1}^i)_{i=1}^N$.

Proposing and selecting the particles according to the dynamics of the latent process, i.e. without making use of the information about the current state provided by the current observation, by letting $R_k \equiv Q$ and $\Psi_k \equiv \mathbf{1}$ for all k , corresponds to the bootstrap particle filter proposed by [Gordon et al. \(1993\)](#).

The algorithm, which was developed gradually by, mainly, [Handschin and Mayne \(1969\)](#), [Gordon et al. \(1993\)](#), and [Pitt and Shephard \(1999\)](#), will be referred to as the *auxiliary particle smoother* (APS). In the setting of a partially observed diffusion process we do not have access to a closed-form expression of the transition density q , which is needed when evaluating the importance weight function Φ_{k+1} . However, the GPE makes it possible to estimate this density without bias via the kernel P . This yields following algorithm, in following referred to as the *GPE-based particle smoother* (GPEPS), in which q in the weighting operation (2.9) is replaced by the Monte Carlo estimate

$$q^\alpha(x, x') \stackrel{\text{def}}{=} \frac{1}{\alpha} \sum_{\ell=1}^{\alpha} V^\ell(x, x'), \quad (2.10)$$

where the $V^\ell(x, x')$'s are drawn independently from $P(x, x', \cdot)$. Denote by

$$\Phi_{k+1}^\alpha(x_{0:k+1}) \stackrel{\text{def}}{=} g(x_{k+1}, Y_{k+1}) \Psi_k^{-1}(x_{0:k}) \frac{q^\alpha(x_k, x_{k+1})}{r_k(x_k, x_{k+1})}, \quad (2.11)$$

the resulting estimated importance weight function. One iteration of the GPEPS is described in detail in the following scheme.

Algorithm 1

(* One iteration of GPEPS *)

Input: $(\xi_{0:k|k}^i, \omega_k^i)_{i=1}^N$, R_k , α

1. **for** $i \leftarrow 1$ **to** N
2. simulate $I_k^i \sim (\omega_k^j \psi_k^j / \sum_{\ell=1}^N \omega_k^\ell \psi_k^\ell)_{j=1}^N$;
3. simulate $\xi_{k+1|k+1}^i \sim R_k(\xi_{k|k}^{I_k^i}, \cdot)$;
4. set $\xi_{0:k+1|k+1}^i \leftarrow (\xi_{0:k|k}^{I_k^i}, \xi_{k+1|k+1}^i)$;
5. simulate $V^{1:\alpha}(\xi_{k:k+1|k+1}^i) \sim P^{\otimes \alpha}(\xi_{k:k+1|k+1}^i, \cdot)$;
6. compute Φ_{k+1}^α via (2.11);
7. set $\omega_{k+1}^i \leftarrow \Phi_{k+1}^\alpha(\xi_{k:k+1|k+1}^i)$;
8. **return** $(\xi_{0:k+1|k+1}^i, \omega_{k+1}^i)_{i=1}^N$.

Here we have used the notations $V^{1:\alpha}(x, x') \stackrel{\text{def}}{=} (V^1(x, x'), \dots, V^\alpha(x, x'))$ and $P^{\otimes \alpha}(x, x', \cdot) \stackrel{\text{def}}{=} P(x, x', \cdot) \otimes \dots \otimes P(x, x', \cdot)$ (α times). Algorithm 1 extends the *random weight auxiliary particle filter* proposed by Fearnhead et al. (2008) to the smoothing mode. Note that we have, in the scheme above, suppressed the dependence of the particles and the particle weights on α from the notation for clarity.

In the selection operation of Step (2), each particle index is drawn from the probability distribution formed by the adjusted weights $(\omega_k^j \psi_k^j / \sum_{\ell=1}^N \omega_k^\ell \psi_k^\ell)_{j=1}^N$. Letting M_k^i denote the number of times that index i was drawn, the selection operation may be alternatively expressed as

$$(M_k^1, \dots, M_k^N) \sim \text{Mult} \left(N, \left(\frac{\omega_k^j \psi_k^j}{\sum_{\ell=1}^N \omega_k^\ell \psi_k^\ell} \right)_{j=1}^N \right). \quad (2.12)$$

There are however many alternative ways of performing selection; e.g., one may set $M_k^i \stackrel{\text{def}}{=} \lfloor N \omega_k^i \psi_k^i / \sum_{\ell=1}^N \omega_k^\ell \psi_k^\ell \rfloor + H_k^i$ with

$$(H_k^1, \dots, H_k^N) \sim \text{Mult} \left(\sum_{i=1}^N \left\langle \frac{N \omega_k^i \psi_k^i}{\sum_{\ell=1}^N \omega_k^\ell \psi_k^\ell} \right\rangle, \left(\frac{\langle N \omega_k^i \psi_k^i / \sum_{\ell=1}^N \omega_k^\ell \psi_k^\ell \rangle}{\sum_{j=1}^N \langle N \omega_k^j \psi_k^j / \sum_{\ell=1}^N \omega_k^\ell \psi_k^\ell \rangle} \right)_{i=1}^N \right), \quad (2.13)$$

where $\lfloor x \rfloor$ denotes the integer part of a real number x and $\langle x \rangle \stackrel{\text{def}}{=} x - \lfloor x \rfloor$. In this selection schedule, which was proposed by Liu and Chen (1995) under the name *deterministic plus residual multinomial resampling*, index i is first copied $\lfloor N \omega_k^i \psi_k^i / \sum_{\ell=1}^N \omega_k^\ell \psi_k^\ell \rfloor$ times; the remaining $\sum_{i=1}^N \langle N \omega_k^i \psi_k^i / \sum_{\ell=1}^N \omega_k^\ell \psi_k^\ell \rangle$ indices are hereafter drawn multinomially with respect to weights proportional to the residuals $(\langle N \omega_k^i \psi_k^i / \sum_{\ell=1}^N \omega_k^\ell \psi_k^\ell \rangle)_{i=1}^N$. All theoretical results obtained in the following will hold for both the selection schedules (2.12) and (2.13). In addition,

our results are easily extended to selection schemes based on *Poisson, binomial, and Bernoulli branching* (see Douc and Moulines, 2008, for a theoretical analysis of these algorithms); however, since the number of drawn indices are random in this case, we omit these results for brevity.

2.2.1. Convergence of the GPEPS

We will describe the convergence, as N tends to infinity, of the self-normalised Monte Carlo approximations formed by weighted particle samples returned by Algorithm 1 using the concept of *consistency* (adopted from Douc and Moulines, 2008) defined in the following. Let $(\Xi, \mathcal{B}(\Xi))$ denote some given state space and $(\xi_{N,i}, \omega_{N,i})_{i=1}^N$ a Ξ -valued particle sample.

Definition 2.1. A weighted sample $(\xi_{N,i}, \omega_{N,i})_{i=1}^N$ is consistent for a probability measure μ and a set $\mathcal{C} \subseteq \mathcal{L}^1(\Xi, \mu)$ if, as $N \rightarrow \infty$,

$$\Omega_N^{-1} \sum_{i=1}^N \omega_{N,i} f(\xi_{N,i}) \xrightarrow{\mathbb{P}} \mu(f), \quad \text{for all } f \in \mathcal{C}, \quad (2.14)$$

and, additionally,

$$\Omega_N^{-1} \max_{1 \leq i \leq N} \omega_{N,i} \xrightarrow{\mathbb{P}} 0. \quad (2.15)$$

The following assumption is mild (in fact, minimal) but essential when establishing consistency of the GPEPS scheme.

(A4) For all $0 \leq k \leq n$, $\Psi_k \in \mathcal{L}^1(\mathbf{X}^{k+1}, \phi_k)$ and $L_k(\cdot, \mathbf{X}) \in \mathcal{L}^1(\mathbf{X}, \phi_{k|k})$.

Proposition 2.1. Assume (A1–4) and that the initial sample $(\xi_0^i, \omega_0^i)_{i=1}^N$ is consistent for $(\phi_0, \mathcal{L}^1(\mathbf{X}, \phi_0))$. Then, for all $1 \leq k \leq n$, each sample $(\xi_{0:k}^i, \omega_k^i)_{i=1}^N$ produced by Algorithm 1 is consistent for $(\phi_k, \mathcal{L}^1(\mathbf{X}^{k+1}, \phi_k))$. The same is true when the multinomial selection schedule (2.12) is replaced by deterministic plus residual multinomial selection (2.13).

The proof of Proposition 2.1 is postponed to Appendix A.1.

2.3. Fixed-lag smoothing

Unfortunately, it has been observed by several authors that using standard SMC methods in the smoothing mode may be unreliable for larger observation sample sizes n , since resampling systematically the particles degenerates the particle paths. Indeed, when $k \ll n$, most (or possibly all) marginal particles $(\xi_{k|n}^i)_{i=1}^N$ will coincide, resulting in a significant Monte Carlo error when estimating any expectation of X_k given $Y_{0:n}$ using the produced particles. Especially, returning to the problem of estimating the intermediate quantity \mathcal{Q}_n in (2.2), for any type of *additive functional* $t(x_{0:n}) \stackrel{\text{def}}{=} \sum_{k=0}^{n-1} s_k(x_{k:k+1})$, $(s_k)_{k=0}^{n-1}$ being a set of

functions (cf. the two terms of (2.2)), we may expect that the estimator

$$(\Omega_n^N)^{-1} \sum_{k=0}^{n-1} \sum_{i=1}^N \omega_n^i s_k(\xi_{k:k+1|n}^i) \quad (2.16)$$

of $\mathbb{E}[t(X_{0:n})|Y_{0:n}]$ is poor when n is large. To compensate for this degeneracy the particle sample size N has to be increased drastically, yielding a computationally inefficient algorithm.

On the other hand, since we may expect that remote observations are only weakly dependent, it should hold that, for a large enough integer Δ_n ,

$$\mathbb{E}[s_k(X_{k:k+1})|Y_{0:n}] \approx \mathbb{E}[s_k(X_{k:k+1})|Y_{0:k(\Delta_n)}] ,$$

where $k(\Delta_n) \stackrel{\text{def}}{=} \min\{k + \Delta_n, n\}$, yielding

$$\mathbb{E}[t(X_{0:n})|Y_{0:n}] = \sum_{k=0}^{n-1} \mathbb{E}[s_k(X_{k:k+1})|Y_{0:n}] \approx \sum_{k=0}^{n-1} \mathbb{E}[s_k(X_{k:k+1})|Y_{0:k(\Delta_n)}] . \quad (2.17)$$

Thus, as long as the approximation (2.17) is relatively precise for a Δ_n which is smaller than the average particle trajectory collapsing time, i.e. most marginal particles $(\xi_{k|k(\Delta_n)}^i)_{i=1}^N$ are different for all k , we should replace (2.16) by the estimator

$$\sum_{k=0}^{n-1} \left(\Omega_{k(\Delta_n)}^N \right)^{-1} \sum_{i=1}^N \omega_{k(\Delta_n)}^i s_k \left(\xi_{k:k+1|k(\Delta_n)}^i \right) . \quad (2.18)$$

The lag-based approximation (2.18) may be computed recursively in a *single sweep* of the data with only limited computer data storage demands, and computing (2.18) is clearly not more computationally demanding than computing (2.16) (having $O(nM)$ complexity); see Olsson et al. (2008) for details. Finally, using (2.18) in conjunction with the kernel \bar{P}_θ for estimating $\log q_\theta$ gives us the following approximation of the intermediate quantity $\mathcal{Q}_n(\theta; \theta')$:

$$\mathcal{Q}_n^N(\theta; \theta') \stackrel{\text{def}}{=} \sum_{k=0}^{n-1} \left(\Omega_{k(\Delta_n)}^{N, \theta'} \right)^{-1} \sum_{i=1}^N \omega_{k(\Delta_n)}^{i, \theta'} s_k^{\bar{\alpha}} \left(\xi_{k:k+1|k(\Delta_n)}^{i, \theta'}; \theta \right) , \quad (2.19)$$

where, for $(x, x') \in \mathbb{X}^2$,

$$s_k^{\bar{\alpha}}(x, x'; \theta) \stackrel{\text{def}}{=} \frac{1}{\bar{\alpha}} \sum_{\ell=1}^{\bar{\alpha}} \bar{V}_\theta^\ell(x, x') + \log g_\theta(x', Y_{k+1})$$

and

$$\bar{V}_\theta^{1:\bar{\alpha}}(x, x') \sim \bar{P}_\theta^{\otimes \bar{\alpha}}(x, x', \cdot) .$$

In (2.19) we have added θ' as an index to the particles as well as the associated weights to indicate that the particle system of the fixed-lag smoother is evolved under the dynamics determined by the initial parameter value.

2.3.1. Convergence of the intermediate quantity

Under weak assumptions on the functions Ψ_k , the kernels L_k and \bar{P} , and the local likelihoods functions $\log g_\theta(\cdot, Y_k)$ one may establish the convergence of the approximate intermediate quantity (2.19). Thus, define, for a given lag Δ_n and parameters (θ, θ') , the bias

$$b_n(\Delta_n, \theta, \theta') \stackrel{\text{def}}{=} \sum_{k=0}^{n-1} \int s_k(x_{k:k+1}, \theta) \phi_{k(\Delta_n)}(dx_{k:k+1}, \theta') - \sum_{k=0}^{n-1} \int s_k(x_{k:k+1}, \theta) \phi_n(dx_{k:k+1}, \theta') \quad (2.20)$$

imposed by the fixed lag. We then have the following result, which is the main result of this section.

Theorem 2.1. *Assume (A1–3). Let $n \geq 0$, $(\theta, \theta') \in \Theta^2$, and $(\Delta_n, \alpha, \bar{\alpha}) \in \mathbb{N}^3$. Suppose that (A4) holds for $\Psi_k(\cdot; \theta')$, $L_k(\cdot; \theta')$, and $\phi_k(\cdot; \theta')$ and that the initial sample $(\xi_0^{i, \theta'}, \omega_0^{i, \theta'})_{i=1}^N$ is consistent for $(\phi_0(\cdot; \theta'), \mathbb{L}^1(\phi_0(\cdot; \theta'), \mathbf{X}))$. Moreover, assume that the mappings $x_{0:k(\Delta_n)} \mapsto \log g_\theta(x_k, Y_k)$, $0 \leq k \leq n$, and $x_{0:k(\Delta_n)} \mapsto \int |v| \bar{P}_\theta(x_k, x_{k+1}, dv)$, $0 \leq k < n$, belong to $\mathbb{L}^1(\phi_{k(\Delta_n)}(\cdot; \theta'), \mathbf{X}^{k(\Delta_n)+1})$. Then, as $N \rightarrow \infty$,*

$$\mathcal{Q}_n^N(\theta, \theta') \xrightarrow{\mathbb{P}} \mathcal{Q}_n(\theta, \theta') + b_n(\Delta_n, \theta, \theta'),$$

where the bias b_n is defined in (2.20).

The proof is given in Appendix A.2.

The bias term b_n , which was treated by Olsson et al. (2008), is controlled by the speed with which the hidden chain $(X_k)_{k \geq 0}$ forgets its initial distribution when evolving *conditionally* on the observations. Indeed, when the state space \mathbf{X} is compact it can be shown (see Olsson et al., 2008, for details) that b_n is $\mathcal{O}(n\rho^{\Delta_n})$, where $0 < \rho < 1$ is the *uniform* (with respect to observation records $Y_{0:n}$ as well as initial distributions χ) mixing coefficient of the conditional chain. From this we deduce that the lag Δ_n should be increased with n at the minimum rate $c \log n$, $c > -1/\log \rho$ in order to keep the bias suppressed. Increasing Δ_n faster eliminates the bias and increases the variance of the approximation; see again Olsson et al. (2008) for a detailed study of these issues. Since a similar forgetting property holds also in the case of a non-compact state space \mathbf{X} (Douc et al., 2009a), the same arguments can be applied for very general models; however, the analysis of the general case is significantly more involved, since the mixing coefficient is neither uniform with respect to observation records nor initial distributions χ in this case.

Remarkably, the convergence result in Theorem 2.1 holds for *any fixed sample sizes* $(\alpha, \bar{\alpha})$. In particular, nothing prevents us from letting $\alpha = \bar{\alpha} = 1$, yielding a computationally very efficient algorithm; this is the choice of Section 3.

2.4. Forward-filtering backward-smoothing

Even though naive SMC implementations generally fail to estimate joint smoothing distributions efficiently, they can, as discussed above, be successfully used for estimating the marginal filter distributions (corresponding to $k = n$ in the discussion of Section 2.3). Nevertheless, any joint smoothing distribution may be expressed in terms of marginal filter distributions via the so-called *forward-filtering backward-smoothing decomposition*. Indeed, for any probability measure η on $(\mathbf{X}, \mathcal{X})$, define the *reverse kernel*

$$\overleftarrow{Q}_\eta(x', A; \theta) \stackrel{\text{def}}{=} \frac{\int_A q_\theta(x, x') \eta(dx)}{\int q_\theta(x, x') \eta(dx)}, \quad (2.21)$$

where $A \in \mathcal{X}$ and $x' \in \mathbf{X}$. The definition (2.21) is valid only when x' belongs to the subset of \mathbf{X} where the denominator is nonzero; outside this set we may let \overleftarrow{Q}_η take arbitrary values. It can now be shown that (see e.g. Cappé et al., 2005, Corollary 3.3.8)

$$\phi_n(A; \theta) = \int \cdots \int_A \phi_{n|n}(dx_n; \theta) \prod_{k=0}^{n-1} \overleftarrow{Q}_{\phi_{k|k}}(x_{k+1}, dx_k; \theta), \quad (2.22)$$

for $A \in \mathcal{X}^{\otimes(n+1)}$. Using the Markovian structure of the decomposition above, a trajectory $X_{0:n}$ can be simulated from $\phi_n(\cdot; \theta)$ by, firstly, computing recursively (via (2.4)) the filter distributions $(\phi_{k|k}(\cdot; \theta))_{k=0}^n$ and, secondly, simulating X_n from $\phi_{n|n}(\cdot; \theta)$ and hereafter, recursively for $k = n-1, n-1, \dots, 0$, X_k from $\overleftarrow{Q}_{\phi_{k|k}}(X_{k+1}, \cdot; \theta)$. This scheme will in the following be referred to as *forward-filtering backward-simulation* (FFBS), and we refer again to Cappé et al. (2005) for a detailed treatment.

In general we lack closed-form expressions of the filter distributions, but may estimate these efficiently using Algorithm 1. Hence, following Doucet et al. (2000), a non-degenerate particle estimate of $\phi_{0:n}(\cdot; \theta)$ can be obtained by replacing, in the decomposition (2.22), $\phi_{n|n}$ by the empirical measure $\phi_{n|n}^N$ and the reverse kernels $\overleftarrow{Q}_{\phi_{k|k}}(x_{k+1}, dx_k; \theta)$ by

$$\overleftarrow{Q}_{\phi_{k|k}^N}(x_{k+1}, dx_k; \theta) = \sum_{i=1}^N \frac{\omega_k^i q_\theta(\xi_{k|k}^i, x_{k+1})}{\sum_{\ell=1}^N \omega_k^\ell q_\theta(\xi_{k|k}^\ell, x_{k+1})} \delta_{\xi_{k|k}^i}(dx_k). \quad (2.23)$$

Note that a draw according to $\overleftarrow{Q}_{\phi_{k|k}^N}(x_{k+1}, \cdot; \theta)$ consists of selecting position $\xi_{k|k}^i$ with probability proportional $\omega_k^i q_\theta(\xi_{k|k}^i, x_{k+1}) / \sum_{\ell=1}^N \omega_k^\ell q_\theta(\xi_{k|k}^\ell, x_{k+1})$. In the case of PODs, a closed-form expression of q_θ is in general missing, and we thus replace each number $q_\theta(\xi_{k|k}^i, x_{k+1})$ by a draw $V_\theta(\xi_{k|k}^i, x_{k+1})$ from the GPE $P_\theta(\xi_{k|k}^i, x_{k+1}, \cdot)$. This gives us the following algorithm for simulating a trajectory $X_{0:n}$ that is approximately distributed according to ϕ_n .

Algorithm 2

(* GPE-based particle FFBS *)

Input: $(R_k)_{k=0}^{n-1}$

1. run Algorithm 1 to obtain $(\phi_{k|k}^N(\cdot; \theta))_{k=0}^n$;
2. simulate $X_n \sim \phi_{n|n}^N(\cdot; \theta)$;
3. **for** $k \leftarrow n-1$ **to** 0
4. **for** $i \leftarrow 1$ **to** N
5. simulate $V_\theta(\xi_{k|k}^i, X_{k+1}) \sim P_\theta(\xi_{k|k}^i, X_{k+1}, \cdot)$;
6. simulate $\iota_k \sim (\omega_k^i V_\theta(\xi_{k|k}^i, X_{k+1}) / \sum_{\ell=1}^N \omega_k^\ell V_\theta(\xi_{k|k}^\ell, X_{k+1}))_{i=1}^N$;
7. set $X_k \leftarrow \xi_{k|k}^{\iota_k}$
8. **return** $X_{0:n} = (X_0, \dots, X_n)$.

Algorithm 2 avoids the problem of degeneracy of the genealogical tree *without* any implicit assumption on geometrical ergodicity of the conditional hidden chain. On the other hand, simulating a single trajectory according to Algorithm 2 involves $\mathcal{O}(N)$ operations, implying an overall computational cost of order $\mathcal{O}(N^2)$ for producing a sample of size N . Recently, Douc et al. (2009b) showed how the overall computational cost of the particle-based FFBS can be reduced to $\mathcal{O}(N)$ by means of accept-reject-methods; however, it is not straightforward to adapt this approach to our framework, since one for general PODs cannot find an upper bound on the transition density of the hidden chain. For models with forgetting properties, Algorithm 2 should be outperformed by the fixed-lag smoother because of the quadratic complexity of the former scheme (see the coming section for examples); the FFBS should thus be seen as a generic and alternative solution in cases of poor mixing.

3. Simulation study

In this section, the proposed methods are illustrated on two simulated examples, consisting of noisy observations of the models treated by Beskos et al. (2006) and Beskos et al. (2008). In both examples we let, for simplicity, the measurement noise variance σ_ϵ be known and set to 0.1 and assume equidistant measurements with $t_{k+1} - t_k = 1$ for all $k \geq 0$. We use consequently $\alpha = \bar{\alpha} = 1$. The approximate intermediate quantity Q_n^N is maximised using the *Nelder-Mead simplex algorithm* as implemented in MATLAB's `fminsearch`-command. In order to obtain convergence of the parameter sequence returned by the Monte Carlo EM-algorithm, it is necessary to decrease, at each iteration, the bias of the particle approximation by increasing the number of particles with the iteration index. We thus follow the recommendations of Fort and Moulines (2003) and increase the particle sample size as the square root of the iteration number, with an initial size of 100 particles. A detailed discussion on the effect of the lag size on the quality of the final parameter estimates is given in Olsson et al. (2008); thus, we do not repeat this discussion here and stick consequently to

the recommendation of increasing the lag logarithmically with the size of the observation record.

3.1. Log-growth model

In the first example we estimate, from simulated data, the parameters of a partially observed version of the *log-growth model* discussed by Beskos et al. (2006). The model is specified by the following system of equations:

$$\begin{aligned} dX_t &= \kappa X_t(1 - X_t/\Lambda) dt + \sigma X_t dW_t, \\ Y_k &= X_{t_k} + \sigma \epsilon_k, \end{aligned} \quad (3.1)$$

where $(\epsilon_k)_{k \geq 0}$ are mutually independent, standard normal-distributed random variables. The noise sequence is supposed to be independent also from W . Applying Itô's formula to the transformation $\tilde{X}_t = \eta(X_t, \sigma)$, with $\eta(x, \sigma) \stackrel{\text{def}}{=} -\log(x)/\sigma$, yields

$$d\tilde{X}_t = \alpha(\tilde{X}_t) dt + dW_t, \quad (3.2)$$

where $\alpha(x) \stackrel{\text{def}}{=} \sigma/2 - \kappa/\sigma + \kappa/(\sigma\Lambda) \exp(-\sigma x)$. Since α is bounded from above, we are only required to simulate the minimum of the Brownian path and let \tilde{W}_α^- be α evaluated at this minimum; see Section B for the meaning of \tilde{W}_α^- . The minimum of the Brownian bridge has a known law, and given the minimum, the bridge can be constructed retrospectively using Bessel bridges (see Beskos et al., 2006). Our aim is to estimate the unknown parameters $\theta \stackrel{\text{def}}{=} (\kappa, \Lambda, \sigma)$ given a record $Y_{0:1000}$ of observations. The observation set was obtained through simulation under the parameters $\theta^* = (0.1, 1000, 0.1)$. When computing the approximate intermediate quantity \mathcal{Q}_n^N , the random weight fixed-lag smoother used the lag $\Delta_n = 40$ and the proposal

$$R_k(x, A) = \frac{1}{\sigma x} \int_A t(\{x' - \kappa x(1 - x/\Lambda)\}/\{\sigma x\}; 4) dx', \quad (3.3)$$

where $t(\cdot; n)$ denotes the density of the student's t -distribution with n degrees of freedom. Further the adjustment multiplier weights are set to 1. The proposal (3.3) is obtained by discretising the hidden dynamics using the Euler scheme. We set $\alpha = \bar{\alpha} = 1$. The EM output is presented in Figure 3.1.

For comparison, the estimation problem of the log-growth model was also solved using the GPE-based particle FFBS in Section 2.4. The setup was the same as for the fixed-lag smoother, but due to the significant higher computational cost of the FFBS scheme (recall Section 2.4) the number of observations was reduced to 100. For the FFBS-based procedure, the GPE needs to be evaluated $N + 1$ times per particle and time step, i.e., once in the forward filtering pass and N times in the backward simulation sweep, compared to only once for the fixed-lag smoother.

The output of the EM learning curves obtained using the GPE-based particle FFBS is presented in Figure 3.1.

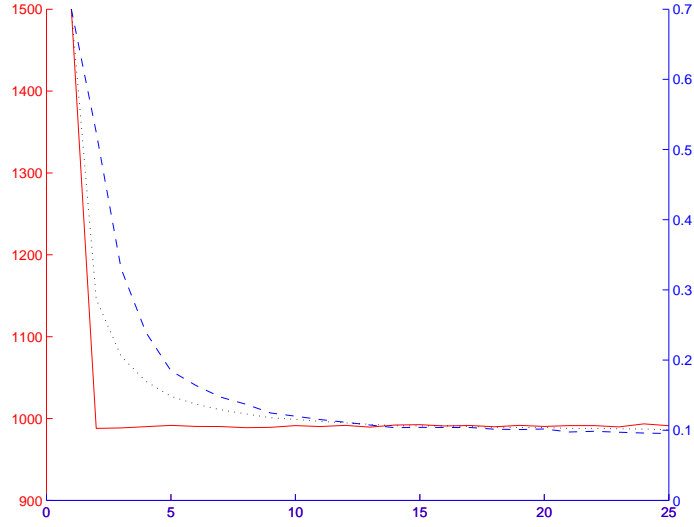


FIG 1. Convergence of Λ (solid, left y-axis), κ (dashed, right y-axis), and σ (dotted, right y-axis) using the fixed-lag smoother with lag 40.

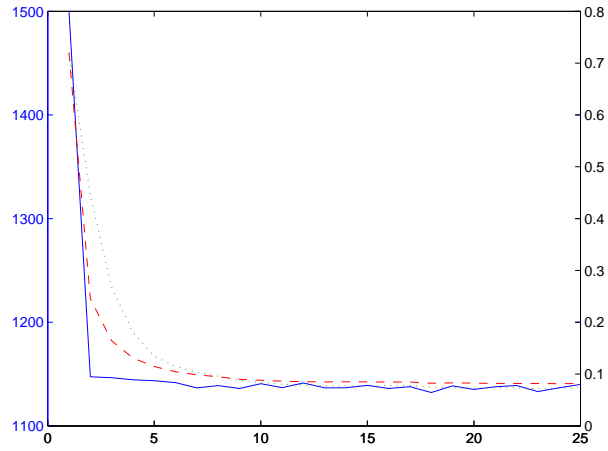


FIG 2. Convergence of Λ (solid, left y-axis), κ (dashed, right y-axis), and σ (dotted, right y-axis) using the GPE-based particle FFBS on 100 observations.

3.2. Genetics diffusion model

In a second example we estimate, again from simulated data, the parameters of a partially observed version of the *genetics diffusion model* presented in Kloeden and Platen (1992) and discussed by Beskos et al. (2008). The model is given by

$$\begin{aligned} dV_t &= (\mu + \nu V_t) dt + \sigma V_t(1 - V_t) dW_t, \\ Y_{t_k} &= V_{t_k} + \sigma_\epsilon \epsilon_k, \end{aligned} \quad (3.4)$$

where the sequence $(\epsilon_k)_{k \geq 0}$ is as in the previous example. Applying Itô's formula to the transformation $\tilde{X}_t = \eta(V_t, \sigma)$, where $\eta(v, \sigma) \stackrel{\text{def}}{=} (\log(v) - \log(1 - v))/\sigma$, allows for using the GPE for estimating the transition density of the latent process. In this case, the drift function α of the transformed process becomes more involved than in the previous example, and it is neither bounded from above nor below. Thus, we have to draw both \tilde{W}_α^- and \tilde{W}_α^+ and a Brownian bridge $(\tilde{W}_s)_{s=0}^t$ such that $\tilde{W}_\alpha^- \leq \alpha(\tilde{W}_s) \leq \tilde{W}_\alpha^+$ for all $0 \leq s \leq t$; see Section B for a justification of this. For this purpose we apply the method proposed in Beskos et al. (2008), which involves sampling first a maximum \tilde{W}_{id}^+ and a minimum \tilde{W}_{id}^- , and then a Brownian bridge such that $\tilde{W}_{\text{id}}^- \leq \tilde{W}_s \leq \tilde{W}_{\text{id}}^+$ for all $0 \leq s \leq t$. Since a linear transformation of a Brownian bridge is still a Brownian bridge, it suffices to consider the case when the path $(\tilde{W}_s)_{s=0}^t$ is conditioned to start and end in zero. Sampling a lower and upper bound can then be done by using rejection sampling in the following way: let $(a_i)_{i \geq 0}$ with $a_0 = 0$ be an increasing sequence and consider the intervals $(-a_i, a_i]$. Since the probability that a Brownian bridge stays in a specific interval $[-K, K]$ has a known expression (having the form of an infinite series), it is possible to calculate the probability that it is contained in $(-a_i, a_i]$ but not in $(-a_{i-1}, a_{i-1}]$; this means that either its maximum is contained in $(a_{i-1}, a_i]$ or its minimum is contained in $(-a_i, -a_{i-1}]$ or both. Thus, we first propose an interval $(a_{i-1}, a_i]$; given this interval, we then propose, with probability 1/2, a maximum conditioned to belong to $(a_{i-1}, a_i]$, otherwise a minimum in $(-a_i, -a_{i-1}]$. Since the distributions of the maximum and minimum are known on closed-form, this is easily done. Next, we propose a Brownian bridge by decomposing around the proposed maximum (minimum) as in the previous example. The resulting path $(\tilde{W}_s)_{s=0}^t$ is accepted, with a probability depending on the path in question, only if it remains in the interval; see Beskos et al. (2008) for details. Finally, we set $\tilde{W}_\alpha^\pm \stackrel{\text{def}}{=} \alpha(\tilde{W}_{\text{id}}^\pm)$.

Again we attempt to estimate the unknown parameters $\theta \stackrel{\text{def}}{=} (\mu, \nu, \sigma)$ given a record $Y_{0:1000}$ of observations obtained through simulation under the parameters $\theta^* = (0.05, 0.1, 1)$. When computing the approximate intermediate quantity \mathcal{Q}_n^N , the random weight fixed-lag smoother used the lag $\Delta_n = 20$. Since the state space $\mathbb{R}(0, 1)$ is compact, we propose the particles by simply drawing uniforms over $(0, 1)$. We set $\alpha = \bar{\alpha} = 1$. The EM output is presented in Figure 3.2.

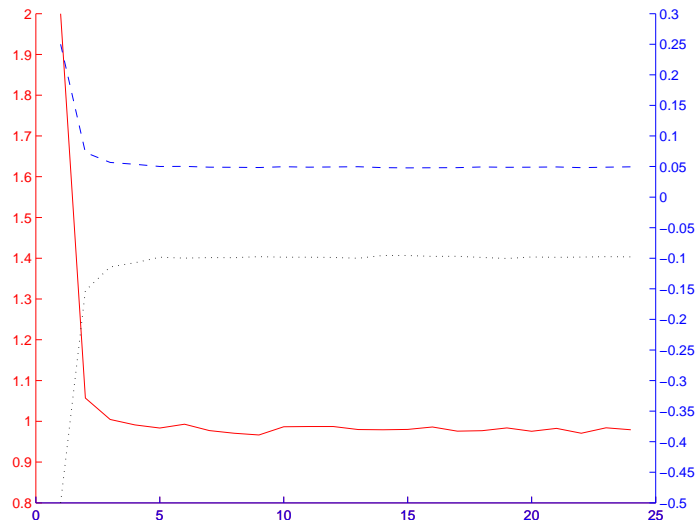


FIG 3. Convergence of σ (dotted, left y-axis), ν (dashed, right y-axis) and μ (dotted, right y-axis).

4. Conclusion

Parameter inference in general discretely and partially observed diffusion processes is an inherently difficult problem due to the lack of closed-form transition densities of the hidden Markov chain. Assuming the possibility of simulating exactly transitions of the latent diffusion process, it is possible to produce pointwise and consistent estimates of the likelihood function using the standard bootstrap particle filter, in which the particles are assigned importance weights determined completely by the known local likelihood function. In such a framework, the likelihood surface can be explored using e.g. grid-based methods (Olsson and Rydén, 2008). Ionides et al. (2009) use the bootstrap particle filter for computing pointwise approximations of the score function and locate the maximum likelihood estimate by means of stochastic approximation. However, simulating exactly transitions of a diffusion process is in general infeasible and we are most often referred to discretisation-based methods such as the Euler scheme, imposing a nontrivially controlled bias of the final parameter estimates. Moreover, mutating blindly, as in the bootstrap particle filter, the particles without incorporating, in the proposal kernel, the information provided by the observations will in general lead to serious degeneracy of the particle weights, especially for models where the observations are informative.

Thus, in the present paper we proposed an alternative, EM-based method for estimating unknown parameters of PODs. The method combines recent approaches for estimating efficiently the joint smoothing distribution in hidden Markov models with recently proposed techniques for estimating, without bias, transition densities of a large class of diffusion processes via GPEs (Beskos et al.,

2008). Interestingly, the GPE provides a way of producing unbiased estimates of the transition densities *simultaneously* for all parameter values; this is critical when carrying through the maximisation-step of the EM-algorithm. For models having forgetting properties, the degeneracy of the particle trajectories can be efficiently avoided by means of fixed-lag smoothing (Kitigawa, 1998; Olsson et al., 2008). The decrease of variance gained by the fixed-lag approximation is obtained at the cost of a bias; the bias is however easily controlled by increasing logarithmically the size of the lag with the size of the observation record, yielding an algorithm of $\mathcal{O}(N)$ computational complexity. We provide a detailed study of the convergence of the GPE-based particle smoother as well as the full intermediate quantity of EM. The results are obtained under, what we believe, minimal assumptions and may, since we analyse separately the GPE-based mutation step (Lemma A.1), be extended to any selection schedule for which consistency has been established in the literature. In this way, our GPEPS convergence results differ significantly from that presented in Fearnhead et al. (2008). In the non-ergodic case, we proposed a method for sampling the joint smoothing distribution which is based on the forward-filtering backward-smoothing decomposition of the same. Basically, the method, which relies on an algorithm proposed by Godsill et al. (2004) and analysed further by Douc et al. (2009b), consists of a forward-filtering pass followed by a backward-simulation pass where trajectories are drawn according to approximations of the backward kernels obtained using the particle filter estimates obtained in the forward pass. During the two passes we replace, when needed, any evaluation of the diffusion process transition density by a draw from the GPE. At the end of the day, we obtain an $\mathcal{O}(N^2)$ algorithm that is significantly more costly than the fixed-lag smoother, but which avoids elegantly the problem of degeneracy of the genealogical tree of the particles. The methods were successfully demonstrated on two examples.

There exist alternative techniques, either Monte Carlo-based (see e.g. Pedersen, 1995) or based on basis expansions (Aït-Sahalia, 2008), for approximating the transition density. Nevertheless, none of these approaches produce unbiased estimates. The former is, while quite general, computationally very demanding and the latter is only valid for very short time intervals (recall that the performance of the GPE is independent of the size of the time grid). Sometimes more direct numerical approaches, such as solving the Fokker-Plank equations or taking the Fourier inverse of the characteristic function of the SDE, are possible; however, these methods often tend to be computationally expensive. Anyway, the theoretical results obtained by us presume only unbiasedness of the transition density estimator, and thus other approximation schemes may be applicable within our framework.

Appendix A: Proofs

The proofs of Proposition 2.1 and Theorem 2.1 rely on recent results on limit theorems for weighted samples obtained by Douc and Moulines (2008). Since we in this section deal exclusively with asymptotic properties of the sample as

the sample size tends to infinity, we let, when not specified differently, the limit notation \rightarrow refer to an *increasing number N of particles* only. In addition, we let also the particles and the associated weights be indexed by N for clearness. The following kernel notation will be useful in the following: Let μ be a measure on $(\Xi, \mathcal{B}(\Xi))$, f a measurable function on $(\tilde{\Xi}, \mathcal{B}(\tilde{\Xi}))$, and K a kernel from $(\Xi, \mathcal{B}(\Xi))$ to $(\tilde{\Xi}, \mathcal{B}(\tilde{\Xi}))$; then we set

$$\mu K(A) \stackrel{\text{def}}{=} \int \mu(d\xi) K(\xi, A)$$

and

$$K(\xi, f) \stackrel{\text{def}}{=} \int f(\tilde{\xi}) K(\xi, d\tilde{\xi}) .$$

The following definition specifies the structure that we want any class of estimand functions to have.

Definition A.1. *A set \mathcal{C} of measurable functions on Ξ is proper if the following holds.*

- (i) \mathcal{C} is a linear space; that is, if f and g belong to \mathcal{C} and $(\alpha, \beta) \in \mathbb{R}^2$, then $\alpha f + \beta g \in \mathcal{C}$;
- (ii) if $g \in \mathcal{C}$ and f is measurable with $|f| \leq |g|$, then $f \in \mathcal{C}$;
- (iii) for all $c \in \mathbb{R}$, the constant function $\xi \mapsto c$ belongs to \mathcal{C} .

We will frequently make use of the following lemma obtained by [Douc and Moulines \(2008\)](#). Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space and $(\mathcal{F}_{N,i})_{i=0}^N$, $N \geq 1$, a triangular array of sub- σ -fields of \mathcal{F} such that $\mathcal{F}_{N,i-1} \subseteq \mathcal{F}_{N,i}$ for all $1 \leq i \leq N$ and $N \geq 1$. In addition, let $(U_{N,i})_{i=1}^N$, $N \geq 1$, be a triangular array of random variables such that each $U_{N,i}$ is $\mathcal{F}_{N,i}$ -measurable.

Theorem A.1 ([Douc and Moulines \(2008\)](#)). *Assume that $\mathbb{E}[|U_{N,j}| | \mathcal{F}_{N,j-1}] < \infty$, \mathbb{P} -a.s., for all $N \geq 1$ and $1 \leq j \leq N$. Suppose that*

- (i) *as $\lambda \rightarrow \infty$,*

$$\sup_{N \geq 1} \mathbb{P} \left(\sum_{j=1}^N \mathbb{E}[|U_{N,j}| | \mathcal{F}_{N,j-1}] \geq \lambda \right) \rightarrow 0 ; \quad (\text{A.1})$$

- (ii) *in addition, for all $\epsilon > 0$,*

$$\sum_{j=1}^N \mathbb{E}[|U_{N,j}|; |U_{N,j}| \geq \epsilon | \mathcal{F}_{N,j-1}] \xrightarrow{\mathbb{P}} 0 \quad (\text{A.2})$$

as $N \rightarrow \infty$. Then

$$\max_{1 \leq i \leq N} \left| \sum_{j=1}^i U_{N,j} - \sum_{j=1}^i \mathbb{E}[U_{N,j} | \mathcal{F}_{N,j-1}] \right| \xrightarrow{\mathbb{P}} 0 .$$

A.1. Proof of Proposition 2.1

Algorithm 1 is conveniently analysed within a more general framework of *random weight mutation* (RWM). Assume that we are given a Ξ -valued, weighted particle sample $(\xi_{N,i}, \omega_{N,i})_{i=1}^N$ which is consistent for some measure ν on $\mathcal{B}(\Xi)$ and let L be a finite transition kernel from $(\Xi, \mathcal{B}(\Xi))$ to $(\tilde{\Xi}, \mathcal{B}(\tilde{\Xi}))$. We wish to transform $(\xi_{N,i}, \omega_{N,i})_{i=1}^N$ into another sample $(\tilde{\xi}_{N,i}, \tilde{\omega}_{N,i})_{i=1}^N$ targeting the measure

$$\mu(A) = \frac{\nu L(A)}{\nu L(\tilde{\Xi})}, \quad A \in \mathcal{B}(\tilde{\Xi}),$$

by means of the RWM operation described below. The input parameters are: a proposal kernel R such that $R(\xi, \cdot)$ dominates $L(\xi, \cdot)$ for all $\xi \in \Xi$, a random weight kernel S from $(\Xi \times \tilde{\Xi}, \mathcal{B}(\Xi \times \tilde{\Xi}))$ to $(\mathbb{R}^+, \mathcal{B}(\mathbb{R}^+))$ targeting dL/dR in the sense that, for all $(\xi, \tilde{\xi}) \in \Xi \times \tilde{\Xi}$,

$$\int v S(\xi, \tilde{\xi}, dv) = \frac{dL(\xi, \cdot)}{dR(\xi, \cdot)}(\tilde{\xi}),$$

and, finally, a Monte Carlo sample size $\alpha \in \mathbb{N}$.

Algorithm 3

(* random weight mutation *)

Input: $(\xi_{N,i}, \omega_{N,i})_{i=1}^N$, R , S , α

1. **for** $i \leftarrow 1$ **to** N
2. **do** simulate $\tilde{\xi}_{N,i} \sim R(\xi_{N,i}, \cdot)$;
3. simulate $V^{1:\alpha}(\xi_{N,i}, \tilde{\xi}_{N,i}) \sim S^{\otimes \alpha}(\xi_{N,i}, \tilde{\xi}_{N,i}, \cdot)$;
4. $\tilde{\omega}_{N,i} \leftarrow \omega_{N,i} \alpha^{-1} \sum_{\ell=1}^{\alpha} V^{\ell}(\xi_{N,i}, \tilde{\xi}_{N,i})$;
5. **return** $(\tilde{\xi}_{N,i}, \tilde{\omega}_{N,i})_{i=1}^N$.

The sample $(\tilde{\xi}_{N,i}, \tilde{\omega}_{N,i})_{i=1}^N$ returned by the algorithm is taken as an approximation of μ . In order to evaluate the quality of this sample, define the set

$$\tilde{\mathcal{C}} \stackrel{\text{def}}{=} \left\{ f \in L^1(\mu, \tilde{\Xi}) : L(\cdot, |f|) \in \mathcal{C} \right\}; \quad (\text{A.3})$$

then the following result stating consistency for weighted samples produced by Algorithm 3 is instrumental when establishing Proposition 2.1.

Lemma A.1. *Assume the weighted sample $(\xi_{N,i}, \omega_{N,i})_{i=1}^N$ is consistent for (ν, \mathcal{C}) and that the function $L(\cdot, \tilde{\Xi})$ belongs to \mathcal{C} . Then the set $\tilde{\mathcal{C}}$ defined in (A.3) and the weighted particle sample $(\tilde{\xi}_{N,i}, \tilde{\omega}_{N,i})_{i=1}^N$ produced by Algorithm 3 are proper resp. $(\mu, \tilde{\mathcal{C}})$ -consistent for any fixed $\alpha \in \mathbb{N}$.*

Proof. Properness of the set $\tilde{\mathcal{C}}$ is straightforwardly established: To check Property (i) in Definition A.1, suppose that f and g belong to $\tilde{\mathcal{C}}$ and let $(\alpha, \beta) \in \mathbb{R}^2$;

then

$$\begin{aligned}
& \iint |\alpha f(\tilde{\xi}) + \beta g(\tilde{\xi})| v S(\cdot, \tilde{\xi}, dv) R(\cdot, d\tilde{\xi}) \\
& \leq |\alpha| \iint |f(\tilde{\xi})| v S(\cdot, \tilde{\xi}, dv) R(\cdot, d\tilde{\xi}) \\
& \quad + |\beta| \iint |g(\tilde{\xi})| v S(\cdot, \tilde{\xi}, dv) R(\cdot, d\tilde{\xi}) \\
& = |\alpha| L(\cdot, |f|) + |\beta| L(\cdot, |g|) ,
\end{aligned}$$

where the function on the right hand side belongs to \mathbb{C} by construction of $\tilde{\mathbb{C}}$ and the fact that \mathbb{C} is a linear space. That the integral on the left hand side belongs to \mathbb{C} is now a consequence of Property (ii) in Definition A.1. Properties (ii) and (iii) are checked in a similar manner.

To establish Condition (2.14) in Definition 2.1 it is enough to show that, for all $f \in \tilde{\mathbb{C}}$,

$$\Omega_N^{-1} \sum_{i=1}^N \tilde{\omega}_{N,i} f(\tilde{\xi}_{N,i}) \xrightarrow{\mathbb{P}} \nu L(f) ; \quad (\text{A.4})$$

indeed, since $\tilde{\mathbb{C}}$ contains the unity mapping $\tilde{\xi} \mapsto 1$ (as $\tilde{\mathbb{C}}$ is proper), (A.4) implies that

$$\Omega_N^{-1} \sum_{i=1}^N \tilde{\omega}_{N,i} \xrightarrow{\mathbb{P}} \nu L(\tilde{\Xi}) , \quad (\text{A.5})$$

from which Condition (2.14) in Definition 2.1 follows by Slutsky's lemma. Thus, we define the triangular array $U_{N,i} \stackrel{\text{def}}{=} \tilde{\omega}_{N,i} f(\tilde{\xi}_{N,i}) / \Omega_N$, $N \geq 1$, $1 \leq i \leq N$, and sub- σ -fields $\mathcal{F}_N \stackrel{\text{def}}{=} \sigma\{(\xi_{N,i}, \omega_{N,i})_{i=1}^N\}$, $N \geq 1$. We then get, by applying the tower property of conditional expectations and the consistency of the ancestor sample,

$$\begin{aligned}
& \sum_{i=1}^N \mathbb{E}[U_{N,i} | \mathcal{F}_N] \\
& = \Omega_N^{-1} \sum_{i=1}^N \omega_{N,i} \mathbb{E} \left[\mathbb{E} \left[\alpha^{-1} \sum_{\ell=1}^{\alpha} V^{\ell}(\xi_{N,i}, \tilde{\xi}_{N,i}) \middle| \tilde{\xi}_{N,i}, \mathcal{F}_N \right] f(\tilde{\xi}_{N,i}) \middle| \mathcal{F}_N \right] \\
& = \Omega_N^{-1} \sum_{i=1}^N \omega_{N,i} \int f(\tilde{\xi}) \int v S(\xi_{N,i}, \tilde{\xi}, dv) R(\xi_{N,i}, d\tilde{\xi}) \\
& = \Omega_N^{-1} \sum_{i=1}^N \omega_{N,i} L(\xi_{N,i}, f) \xrightarrow{\mathbb{P}} \nu L(f) ,
\end{aligned}$$

since $L(\cdot, f) \leq L(\cdot, |f|) \in \mathbb{C}$. To show that $\sum_{i=1}^N U_{N,i}$ tends to $\sum_{i=1}^N \mathbb{E}[U_{N,i} | \mathcal{F}_N]$ in probability, implying (A.4), we apply Theorem A.1. In order to establish

the first condition of that theorem we reuse the arguments above and use that $L(\cdot, |f|) \in \mathbb{C}$, yielding the limit

$$\sum_{i=1}^N \mathbb{E}[|U_{N,i}| | \mathcal{F}_N] \xrightarrow{\mathbb{P}} \nu L(|f|) .$$

Now, since convergence in probability implies tightness, we conclude that Condition (i) in Theorem A.1 is fulfilled.

To verify (ii), define, for some $\epsilon > 0$, $A_N \stackrel{\text{def}}{=} \sum_{i=1}^N \mathbb{E}[|U_{N,i}|; |U_{N,i}| \geq \epsilon | \mathcal{F}_N]$. Since, as the ancestor sample is assumed to be consistent, $\max_{1 \leq i \leq N} \omega_{N,i}/\Omega_N$ vanishes in probability as N tends to infinity, the same holds for the product $A_N \mathbb{1}\{C \max_{1 \leq i \leq N} \omega_{N,i} > \epsilon \Omega_N\}$, where $C > 0$ is an arbitrary constant. On the other hand,

$$\begin{aligned} & A_N \mathbb{1}\left\{C \max_{1 \leq i \leq N} \omega_{N,i} \leq \epsilon \Omega_N\right\} \\ & \leq \sum_{i=1}^N \mathbb{E}\left[|U_{N,i}|; |f(\tilde{\xi}_{N,i})| \sum_{\ell=1}^{\alpha} V^{\ell}(\xi_{N,i}, \tilde{\xi}_{N,i}) \geq \alpha C \middle| \mathcal{F}_N\right] \\ & = \Omega_N^{-1} \sum_{i=1}^N \omega_{N,i} \int |f(\tilde{\xi})| \int_{|f(\tilde{\xi})| \sum_{\ell=1}^{\alpha} v_{\ell} \geq \alpha C} v_1 S^{\otimes \alpha}(\xi_{N,i}, \tilde{\xi}, dv_{1:\alpha}) R(\xi_{N,i}, d\tilde{\xi}) . \end{aligned}$$

Now, since, for all $\xi \in \Xi$,

$$\int |f(\tilde{\xi})| \int_{|f(\tilde{\xi})| \sum_{\ell=1}^{\alpha} v_{\ell} \geq \alpha C} v_1 S^{\otimes \alpha}(\xi, \tilde{\xi}, dv_{1:\alpha}) R(\xi, d\tilde{\xi}) \leq L(\xi, |f|) ,$$

where $L(\cdot, |f|) \in \mathbb{C}$, we conclude, using Property (ii) of Definition A.1, that the mapping

$$\xi \mapsto \int |f(\tilde{\xi})| \int_{|f(\tilde{\xi})| \sum_{\ell=1}^{\alpha} v_{\ell} \geq \alpha C} v_1 S^{\otimes \alpha}(\xi, \tilde{\xi}, dv_{1:\alpha}) R(\xi, d\tilde{\xi})$$

on Ξ belongs to \mathbb{C} as well. Thus, consistency of the ancestor sample implies that

$$\begin{aligned} & \sum_{i=1}^N \mathbb{E}\left[|U_{N,i}|; |f(\tilde{\xi}_{N,i})| \sum_{\ell=1}^{\alpha} V^{\ell}(\xi_{N,i}, \tilde{\xi}_{N,i}) \geq \alpha C \middle| \mathcal{F}_N\right] \\ & \xrightarrow{\mathbb{P}} \iint |f(\tilde{\xi})| \int_{|f(\tilde{\xi})| \sum_{\ell=1}^{\alpha} v_{\ell} \geq \alpha C} v_1 S^{\otimes \alpha}(\xi, \tilde{\xi}, dv_{1:\alpha}) R(\xi, d\tilde{\xi}) \nu(\xi) . \quad (\text{A.6}) \end{aligned}$$

In addition, since the constant C may be chosen arbitrarily large, the limit (A.6) can be made arbitrarily small by the dominated convergence theorem. We hence conclude that A_N tends to zero in probability as N tends to infinity. This establishes (A.4).

In order to establish (2.15) it is, by Slutsky's theorem and (A.5), enough to prove that

$$\Omega_N^{-1} \max_{1 \leq i \leq N} \tilde{\omega}_{N,i} \xrightarrow{\mathbb{P}} 0. \quad (\text{A.7})$$

Thus, take again a constant $C > 0$ and write

$$\begin{aligned} \Omega_N^{-1} \max_{1 \leq i \leq N} \tilde{\omega}_{N,i} \mathbb{1} \left\{ \sum_{\ell=1}^{\alpha} V^{\ell}(\xi_{N,i}, \tilde{\xi}_{N,i}) \geq \alpha C \right\} \\ \leq \Omega_N^{-1} \sum_{i=1}^N \tilde{\omega}_{N,i} \mathbb{1} \left\{ \sum_{\ell=1}^{\alpha} V^{\ell}(\xi_{N,i}, \tilde{\xi}_{N,i}) \geq \alpha C \right\}. \end{aligned} \quad (\text{A.8})$$

To prove that the right hand side of (A.8) converges, we introduce the triangular array $U_{N,i} \stackrel{\text{def}}{=} \tilde{\omega}_{N,i} \mathbb{1} \{ \sum_{\ell=1}^{\alpha} V^{\ell}(\xi_{N,i}, \tilde{\xi}_{N,i}) \geq \alpha C \} / \Omega_N$, $N \geq 1$, $1 \leq i \leq N$, and let the sub- σ -fields \mathcal{F}_N , $N \geq 1$, be defined as above. Next, we use again Theorem A.1. To verify the first condition, take conditional expectation with respect to \mathcal{F}_N and reuse (A.6) with f being the unity function; this yields

$$\sum_{i=1}^N \mathbb{E}[U_{N,i} | \mathcal{F}_N] \xrightarrow{\mathbb{P}} \iiint_{\sum_{\ell=1}^{\alpha} v_{\ell} \geq \alpha C} v_1 S^{\otimes \alpha}(\xi, \tilde{\xi}, dv_{1:\alpha}) R(\xi, d\tilde{\xi}) \nu(d\xi),$$

implying (i). To verify (ii), take an $\epsilon > 0$ and define $A_N \stackrel{\text{def}}{=} \sum_{i=1}^N \mathbb{E}[|U_{N,i}|; |U_{N,i}| \geq \epsilon | \mathcal{F}_N]$. Then

$$A_N = \Omega_N^{-1} \sum_{i=1}^N \omega_{N,i} \mathbb{E} \left[V^1(\xi_{N,i}, \tilde{\xi}_{N,i}); \tilde{\omega}_{N,i} \geq \epsilon \Omega_N, \sum_{\ell=1}^{\alpha} V^{\ell}(\xi_{N,i}, \tilde{\xi}_{N,i}) \geq \alpha C \middle| \mathcal{F}_N \right],$$

implying that, for an arbitrary constant $C' > 0$, following the lines of (A.6),

$$\begin{aligned} A_N \mathbb{1} \left\{ C' \max_{1 \leq i \leq N} \omega_{N,i} \leq \epsilon \Omega_N \right\} \\ \leq \Omega_N^{-1} \sum_{i=1}^N \omega_{N,i} \mathbb{E} \left[V^1(\xi_{N,i}, \tilde{\xi}_{N,i}); \sum_{\ell=1}^{\alpha} V^{\ell}(\xi_{N,i}, \tilde{\xi}_{N,i}) \geq \alpha(C \vee C') \middle| \mathcal{F}_N \right] \\ \xrightarrow{\mathbb{P}} \iiint_{\sum_{\ell=1}^{\alpha} v_{\ell} \geq \alpha(C \vee C')} v_1 S^{\otimes \alpha}(\xi, \tilde{\xi}, dv_{1:\alpha}) R(\xi, d\tilde{\xi}) \nu(d\xi). \end{aligned} \quad (\text{A.9})$$

On the other hand,

$$\Omega_N^{-1} \max_{1 \leq i \leq N} \tilde{\omega}_{N,i} \mathbb{1} \left\{ \sum_{\ell=1}^{\alpha} V^{\ell}(\xi_{N,i}, \tilde{\xi}_{N,i}) < \alpha C \right\} \leq C \Omega_N^{-1} \max_{1 \leq i \leq N} \omega_{N,i} \xrightarrow{\mathbb{P}} 0.$$

Thus, since the limit (A.9) can be made arbitrarily small by increasing C' , we conclude that A_N tends to zero as N tends to infinity. This in turn implies that the upper bound in (A.8) tends to

$$\iiint_{\sum_{\ell=1}^{\alpha} v_{\ell} \geq \alpha C} v_1 S^{\otimes \alpha}(\xi, \tilde{\xi}, dv_{1:\alpha}) R(\xi, d\tilde{\xi}) \nu(d\xi). \quad (\text{A.10})$$

Finally, we complete the proof by noting that (A.10) can be made arbitrarily small by increasing C . \square

We now use Lemma A.1 to prove consistency of Monte Carlo estimates produced by the GPEPS. For this purpose, let $\bar{\xi}_{0:k|k}^i \stackrel{\text{def}}{=} \xi_{0:k|k}^{I_k^i}$, $1 \leq i \leq N$, denote the selected particles obtained in Step (2) of Algorithm 1. Consequently, the sample $(\bar{\xi}_{0:k|k}^i)_{i=1}^N$ is obtained by resampling the ancestor particles $(\xi_{0:k|k}^i)_{i=1}^N$ multinomially with respect to the normalised adjusted weights $(\omega_k^j \psi_k^j / \sum_{\ell=1}^N \omega_k^\ell \psi_k^\ell)_{j=1}^N$. This operation will in the following be referred to as *selection*. Using this notation and terminology it is now possible to describe one iteration of the GPEPS by the following three transformations:

$$\begin{aligned} (\xi_{0:k|k}^i, \omega_k^i)_{i=1}^N &\xrightarrow{\text{I: Weighting}} (\xi_{0:k|k}^i, \psi_k^i \omega_k^i)_{i=1}^N \rightarrow \\ &\xrightarrow{\text{II: Selection}} (\bar{\xi}_{0:k|k}^i, 1)_{i=1}^N \xrightarrow{\text{III: Mutation}} (\xi_{0:k+1|k+1}^i, \omega_{k+1}^i)_{i=1}^N. \end{aligned}$$

Here the third operation refers to the random weight mutation procedure described in Algorithm 3.

To prove Proposition 2.1 we proceed by induction and assume that $(\xi_{0:k|k}^i, \omega_k^i)_{i=1}^N$ is consistent for $(\phi_k, \mathbb{L}^1(\mathbf{X}^{k+1}, \phi_k))$. Next, we show how consistency is preserved through one iteration of the algorithm by analysing separately Steps (I–III).

Step I. Define the modulated smoothing measure

$$\phi_k \langle \Psi_k \rangle (A) \stackrel{\text{def}}{=} \frac{\phi_k(\Psi_k \mathbb{1}_A)}{\phi_k(\Psi_k)}, \quad A \in \mathcal{X}^{\otimes(n+1)};$$

then the weighting operation in Step I can be viewed as a transformation according Algorithm 3 with $\Xi = \mathbf{X}^{n+1}$, $\tilde{\Xi} = \mathbf{X}^{n+1}$, and

$$\begin{cases} \nu = \phi_k, \\ \mu = \phi_k \langle \Psi_k \rangle, \\ R(x_{0:k}, A) = \delta_{x_{0:k}}(A), \\ L(x_{0:k}, A) = \Psi_k(x_{0:k}) \delta_{x_{0:k}}(A), \\ S(x_{0:k}, x'_{0:k}, A) = \delta_{\Psi_k(x'_{0:k})}(A). \end{cases}$$

Thus, by applying Lemma A.1 we conclude that $(\xi_{0:k|k}^i, \psi_k^i \omega_k^i)_{i=1}^N$ is consistent for $\phi_k \langle \Psi_k \rangle$ and the (proper) set

$$\{f \in \mathbb{L}^1(\phi_k \langle \Psi_k \rangle, \mathbf{X}^{n+1}) : \Psi_k |f| \in \mathbb{L}^1(\phi_k, \mathbf{X}^{n+1})\} = \mathbb{L}^1(\phi_k \langle \Psi_k \rangle, \mathbf{X}^{n+1}).$$

Step II. Applying Theorem 3 in Douc and Moulines (2008) gives immediately that $(\bar{\xi}_{0:k|k}^i, 1)_{i=1}^N$ is consistent for $[\phi_k \langle \Psi_k \rangle, \mathbb{L}^1(\phi_k \langle \Psi_k \rangle, \mathbf{X}^{n+1})]$ for both the selection schedules (2.12) and (2.13).

Step III. Also the third step is handled using Lemma A.1. In this case, we set $\Xi = \mathbf{X}^{n+1}$, $\tilde{\Xi} = \mathbf{X}^{n+2}$, and

$$\begin{cases} \nu = \phi_k \langle \Psi_k \rangle, \\ \mu = \phi_{k+1}, \\ R(x_{0:k}, A) = \int_A \delta_{x_{0:k}}(dx'_{0:k}) R_k(x'_k, dx'_{k+1}), \\ L(x_{0:k}, A) = \int_A \Phi_k(x'_{0:k+1}) \delta_{x_{0:k}}(dx'_{0:k}) R_k(x'_k, dx'_{k+1}), \\ S(x_{0:k}, x'_{0:k+1}, A) \\ \quad = \int \mathbb{1}_A \{vg(x'_{k+1}, Y_{k+1}) / [\Psi_k(x'_{0:k}) r_k(x'_k, x'_{k+1})]\} P(x'_k, x'_{k+1}, dv), \end{cases}$$

where P is the GPE described in Section 2.1 (and in more detail in Appendix B). Thus, using Lemma A.1 yields that $(\xi_{0:k+1|k+1}^i, \omega_{k+1}^i)_{i=1}^N$ is consistent for ϕ_{k+1} and the set

$$\{f \in \mathbf{L}^1(\phi_{k+1}, \mathbf{X}^{k+2}) : L(\cdot, |f|) \in \mathbf{L}^1(\phi_k \langle \Psi_k \rangle, \mathbf{X}^{n+1})\} = \mathbf{L}^1(\phi_{k+1}, \mathbf{X}^{k+2}).$$

Finally, we complete the proof by noting that the induction hypothesis is fulfilled for $k = 0$ by assumption.

A.2. Proof of Theorem 2.1

Decompose the error according to

$$\begin{aligned} \mathcal{Q}_n^N(\theta, \theta') - \mathcal{Q}_n(\theta, \theta') &= \sum_{k=0}^{n-1} \left[\left(\Omega_{k(\Delta_n)}^{N, \theta'} \right)^{-1} \sum_{i=1}^N \omega_{k(\Delta_n)}^{i, \theta'} s_{\bar{k}}^{i, \theta'} \left(\xi_{k:k+1|k(\Delta_n)}^{i, \theta'}, \theta \right) \right. \\ &\quad \left. - \int s_k(x_{k:k+1}; \theta) \phi_{k(\Delta_n)}(dx_{k:k+1}; \theta') \right] \\ &\quad + b_n(\Delta_n, \theta, \theta'), \quad (\text{A.11}) \end{aligned}$$

where the bracket terms are errors originating from the GPEPS and the second term b_n , defined in (2.20), is the cost of introducing the fixed lag. By combining Proposition 2.1 with Slutsky's theorem we conclude that

$$\begin{aligned} \sum_{k=0}^n \left(\Omega_{k(\Delta_n)}^{N, \theta'} \right)^{-1} \sum_{i=1}^N \omega_{k(\Delta_n)}^{i, \theta'} \log g_\theta \left(\xi_{k|k(\Delta_n)}^{i, \theta'}, Y_k \right) \\ \xrightarrow{\mathbb{P}} \sum_{k=0}^n \int \log g_\theta(x_k, Y_k) \phi_{k(\Delta_n)}(dx_k; \theta'), \quad (\text{A.12}) \end{aligned}$$

as $x_{0:k(\Delta_n)} \mapsto \log g_\theta(x_k, Y_k)$ belongs to $\mathbf{L}^1(\phi_{k(\Delta_n)}(\cdot; \theta'), \mathbf{X}^{k(\Delta_n)+1})$ by assumption. Thus, the second term of the intermediate quantity estimator (2.19) is

consistent. In order to establish consistency of the complete estimator it remains to prove that

$$\begin{aligned} \sum_{k=0}^{n-1} \left(\bar{\alpha} \Omega_{k(\Delta_n)}^{N, \theta'} \right)^{-1} \sum_{i=1}^N \omega_{k(\Delta_n)}^{i, \theta'} \sum_{\ell=1}^{\bar{\alpha}} \bar{V}_{\theta}^{\ell} \left(\xi_{k:k+1|k(\Delta_n)}^{i, \theta'} \right) \\ \xrightarrow{\mathbb{P}} \sum_{k=0}^{n-1} \int \log q_{\theta}(x_k, x_{k+1}) \phi_{k(\Delta_n)}(dx_{k:k+1}; \theta') . \end{aligned} \quad (\text{A.13})$$

To do this, we define $\bar{U}_{N,i} \stackrel{\text{def}}{=} \omega_{k(\Delta_n)}^{i, \theta'} \sum_{\ell=1}^{\bar{\alpha}} \bar{V}_{\theta}^{\ell}(\xi_{k:k+1|k(\Delta_n)}^{i, \theta'}) / \bar{\alpha} \Omega_{k(\Delta_n)}^{N, \theta'}$ and $\bar{\mathcal{F}}_N \stackrel{\text{def}}{=} \sigma\{(\xi_{0:k(\Delta_n)|k(\Delta_n)}^{i, \theta'}, \omega_{k(\Delta_n)}^{i, \theta'})_{i=1}^N\}$ and appeal to Theorem A.1 and Proposition 2.1. Since $\log q_{\theta}(x_k, x_{k+1}) \leq \int |v| \bar{P}_{\theta}(x_k, x_{k+1}, dv)$ for all $x_{k:k+1} \in \mathbb{X}^2$, the mapping $x_{0:k(\Delta_n)} \mapsto \log q_{\theta}(x_k, x_{k+1})$ belongs to $L^1(\phi_{k(\Delta_n)}(\cdot; \theta'), \mathbb{X}^{k(\Delta_n)+1})$. Hence,

$$\begin{aligned} \sum_{i=1}^N \mathbb{E} [\bar{U}_{N,i} | \bar{\mathcal{F}}_N] &= \left(\Omega_{k(\Delta_n)}^{N, \theta'} \right)^{-1} \sum_{i=1}^N \omega_{k(\Delta_n)}^{i, \theta'} \log q_{\theta} \left(\xi_{k:k+1|k(\Delta_n)}^{i, \theta'} \right) \\ &\xrightarrow{\mathbb{P}} \int \log q_{\theta}(x_k, x_{k+1}) \phi_{k(\Delta_n)}(dx_{k:k+1}; \theta') , \end{aligned} \quad (\text{A.14})$$

from which we conclude that (A.13) may be established by verifying the two assumptions of Theorem A.1. Following (A.14) and using again that $x_{0:k(\Delta_n)} \mapsto \int |v| \bar{P}_{\theta}(x_k, x_{k+1}, dv)$ belongs to $L^1(\phi_{k(\Delta_n)}(\cdot; \theta'), \mathbb{X}^{k(\Delta_n)+1})$ by assumption, we conclude that

$$\sum_{i=1}^N \mathbb{E} [|\bar{U}_{N,i}| | \bar{\mathcal{F}}_N] \xrightarrow{\mathbb{P}} \iint |v| \bar{P}_{\theta}(x_k, x_{k+1}, dv) \phi_{k(\Delta_n)}(dx_{k:k+1}; \theta') ,$$

which verifies Assumption (i) (by tightness of sequences converging in probability). To verify (ii), let $\epsilon > 0$ and set $\bar{A}_N \stackrel{\text{def}}{=} \sum_{i=1}^N \mathbb{E}[|\bar{U}_{N,i}|; |\bar{U}_{N,i}| \geq \epsilon | \bar{\mathcal{F}}_N]$. Then, for any constant $C > 0$, by consistency of the particle sample,

$$\bar{A}_N \mathbb{1} \left\{ C \max_{1 \leq i \leq N} \omega_{k(\Delta_n)}^{i, \theta'} > \epsilon \Omega_{k(\Delta_n)}^{N, \theta'} \right\} \xrightarrow{\mathbb{P}} 0 . \quad (\text{A.15})$$

On the other hand,

$$\begin{aligned} \bar{A}_N \mathbb{1} \left\{ C \max_{1 \leq i \leq N} \omega_{k(\Delta_n)}^{i, \theta'} \leq \epsilon \Omega_{k(\Delta_n)}^{N, \theta'} \right\} \\ \leq \sum_{i=1}^N \mathbb{E} \left[|\bar{U}_{N,i}|; \left| \sum_{\ell=1}^{\bar{\alpha}} \bar{V}_{\theta}^{\ell} \left(\xi_{k:k+1|k(\Delta_n)}^{i, \theta'} \right) \right| \geq C \bar{\alpha} \middle| \bar{\mathcal{F}}_N \right] \\ \leq \left(\Omega_{k(\Delta_n)}^{N, \theta'} \right)^{-1} \sum_{i=1}^N \omega_{k(\Delta_n)}^{i, \theta'} \int_{|\sum_{\ell=1}^{\bar{\alpha}} v_{\ell}| \geq C \bar{\alpha}} |v_1| \bar{P}_{\theta}^{\otimes \bar{\alpha}}(x_k, x_{k+1}, dv_{1:\bar{\alpha}}) . \end{aligned}$$

Now, since, for all $x_{k:k+1} \in \mathbf{X}^2$,

$$\int_{|\sum_{\ell=1}^{\bar{\alpha}} v_{\ell}| \geq C\bar{\alpha}} |v_1| \bar{P}_{\theta}^{\otimes \bar{\alpha}}(x_k, x_{k+1}, dv_{1:\bar{\alpha}}) \leq \int |v| \bar{P}_{\theta}(x_k, x_{k+1}, dv) ,$$

we get, using Proposition 2.1,

$$\begin{aligned} & \left(\Omega_{k(\Delta_n)}^{N, \theta'} \right)^{-1} \sum_{i=1}^N \omega_{k(\Delta_n)}^{i, \theta'} \int_{|\sum_{\ell=1}^{\bar{\alpha}} v_{\ell}| \geq C\bar{\alpha}} |v_1| \bar{P}_{\theta}^{\otimes \bar{\alpha}}(x_k, x_{k+1}, dv_{1:\bar{\alpha}}) \\ & \xrightarrow{\mathbb{P}} \iint_{|\sum_{\ell=1}^{\bar{\alpha}} v_{\ell}| \geq C\bar{\alpha}} |v_1| \bar{P}_{\theta}^{\otimes \bar{\alpha}}(x_k, x_{k+1}, dv_{1:\bar{\alpha}}) \phi_{k(\Delta_n)}(dx_{k:k+1}; \theta') . \end{aligned} \quad (\text{A.16})$$

We now note that the limit in (A.16) can be made arbitrarily small by increasing C . This verifies condition (ii) in Theorem A.1, which completes the proof of (A.13). Finally, combining (A.13) with (A.12) completes the proof of Theorem 2.1.

Appendix B: More on the GPE

The outline of this section follows Beskos et al. (2006) and Fearnhead et al. (2008), and we limit our scope to the one-dimensional case; multivariate extensions are treated by Beskos et al. (2008). Let $(C[0, t], \mathcal{C}[0, t])$ be the measurable space of continuous functions on $[0, t]$ and denote by $\mathfrak{B}x_{\theta}$ the law of \tilde{X} on $(C[0, t], \mathcal{C}[0, t])$ for the initial condition $\tilde{X}_0 = W_0 = x$. Also, let $\mathbb{W}^{(t, x, x')}$ be the law, on the same space, of the Brownian bridge process $\tilde{W} = (\tilde{W}_s)_{0 \leq s \leq t}$ starting in x at time zero and ending in x' at time t . Similarly, denote by $\mathfrak{B}t, x, x'_{\theta}$ the law of the *diffusion bridge* obtained when \tilde{X} is conditioned to start at $\tilde{X}_0 = W_0 = x$ and to finish at $\tilde{X}_t = x'$. Recall the definition (2.1) of $\alpha(\cdot, \theta)$ and let

$$A(u, \theta) \stackrel{\text{def}}{=} \int^u \alpha(v, \theta) dv$$

be any antiderivative of $\alpha(\cdot, \theta)$. The role of Assumptions (A1–A3) is to guarantee that $\mathfrak{B}t, x, x'_{\theta}$ is absolutely continuous with respect to $\mathbb{W}^{(t, x, x')}$ with Radon-Nikodym derivative

$$\begin{aligned} & \frac{d\mathfrak{B}x, x', t_{\theta}}{d\mathbb{W}^{(t, x, x')}}(w) \\ & = \frac{\mathcal{N}_t(x' - x)}{\tilde{q}_{\theta}(x, x', t)} \exp \left(A(x', \theta) - A(x, \theta) - \frac{1}{2} \int_0^t (\alpha^2 + \alpha')(w_s, \theta) ds \right) , \end{aligned} \quad (\text{B.1})$$

where $w \in C[0, t]$ and \mathcal{N}_t denotes the density function of the zero mean normal distribution with variance t . Now, define, for $u \in \mathbb{R}$, the *drift functional*

$$\phi(u, \theta) \stackrel{\text{def}}{=} \frac{\alpha^2(u, \theta) + \alpha'(u, \theta)}{2} - l(\theta) ,$$

where $l(\theta)$ is the lower bound given in Assumption (A3). The transition density \tilde{q}_θ can, using (B.1), be expressed as

$$\begin{aligned} \tilde{q}_\theta(x, x', t) &= \mathcal{N}_t(x' - x) \exp(A(x', \theta) - A(x, \theta) - l(\theta)t) \\ &\quad \times \int \exp\left(-\int_0^t \phi(w_s, \theta) ds\right) \mathbb{W}^{(t, x, x')}(dw), \end{aligned}$$

Accordingly, we wish to calculate expectations of the form

$$\int \exp\left(-\int_0^t f(w_s) ds\right) \mathbb{W}^{(t, x, x')}(dw). \quad (\text{B.2})$$

Now assume that it is possible to simulate simultaneously a pair $(\tilde{W}_f^-, \tilde{W}_f^+)$ of random variables and a trajectory $(\tilde{W}_s)_{s=0}^t$ such that

$$\tilde{W}_f^- \leq f(\tilde{W}_s) \leq \tilde{W}_f^+, \quad \text{for all } s \in [0, t];$$

in practice this will most often be carried through by first simulating a maximum and a minimum of the Brownian bridge process \tilde{W} and hereafter interpolating, using Bessel bridges, the rest of the bridge conditionally on these. Let κ be a discrete random variable having, conditionally on \tilde{W}_f^\pm , probability distribution $p_t(\cdot | \tilde{W}_f^\pm)$. Then it is easily established that the GPE

$$\exp(-\tilde{W}_f^+ t) \frac{t^\kappa}{\kappa! p_t(\kappa | \tilde{W}_f^\pm)} \prod_{\ell=1}^{\kappa} [\tilde{W}_f^+ - f(\tilde{W}_{\psi_\ell})]$$

(associated with p_t) is an unbiased estimator of (B.2). Here $(\psi_\ell)_{\ell \geq 1}$ are mutually independent variables that are uniformly distributed over $[0, t]$ and independent of \mathcal{F}_t . Note that the distribution p_t can be chosen freely, yielding a *whole class* of GPEs, and an optimal choice is discussed by Fearnhead et al. (2008). In all applications considered in this paper we will use let κ be Poisson-distributed.

Using the Girsanov theorem, it can be shown that

$$\begin{aligned} \log \tilde{q}_t(x, x') &= -\frac{1}{2} \log(2\pi t) - \frac{(x' - x)^2}{2t} \\ &\quad + A(x', \theta) - A(x, \theta) - l(\theta)t - \int \left(\int_0^t \phi(w_s, \theta) ds \right) \mathbb{B}x, x', t(dw), \end{aligned} \quad (\text{B.3})$$

Since the right hand side of (B.1) can be bounded from above and below, a rejection sampler producing samples from the diffusion bridge can be constructed. This is possible as the right hand side of (B.1) is proportional to the probability that a marked Poisson process on $[0, t] \times [0, 1]$ with intensity $r \stackrel{\text{def}}{=} \sup_x \{\phi(x); \tilde{W}_\phi^- < x < \tilde{W}_\phi^+\}$ is below the graph $s \mapsto \phi(\tilde{W}_s; \theta)/r$. However, while observing the path for all s is impossible, a finite construction can be devised by sampling the Brownian bridge at points specified by the marked Poisson process; we refer to Beskos et al. (2006) for details. The algorithm is described by the following.

Algorithm 4

(* Sampling a skeleton of a diffusion bridge *)

1. simulate an outcome $(\chi_\ell, \psi_\ell)_{\ell=1}^\kappa$ of the marked Poisson process with intensity r and $\kappa \sim \text{Po}(r)$;
2. conditional on \tilde{W}_ϕ^\pm , simulate $(\tilde{W}_{\chi_\ell})_{\ell=1}^\kappa$;
3. **if** $\phi(\tilde{W}_{\chi_\ell})/r < \psi_\ell$
4. **then return** $(\tilde{W}_{\chi_\ell})_{\ell=1}^\kappa$
5. **else** go to (1)

By interpolating the returned skeleton $(\tilde{W}_{\chi_\ell})_{\ell=1}^\kappa$, samples \tilde{W}_u , with $(\tilde{W}_s)_{s=0}^t \sim \mathbb{B}x, x', t$, can be obtained for any $0 \leq u \leq t$. Given samples from the diffusion bridge, an unbiased estimator of (B.3) can be straightforwardly constructed in the following way. Let $\psi \sim \text{Unif}(0, t)$ be independent of \mathcal{F}_t . Then $-t\phi(\tilde{W}_\psi, \theta)$ is an unbiased estimator of $\int (\int_0^t \phi(w_s, \theta) ds) \mathbb{B}x, x', t(dw)$ since

$$\begin{aligned} \mathbb{E} [t\phi(\tilde{W}_\psi, \theta)] &= \mathbb{E} \left[\mathbb{E} [t\phi(\tilde{W}_\psi, \theta) \mid \mathcal{F}_t] \right] \\ &= \mathbb{E} \int_0^t \phi(\tilde{W}_s, \theta) ds = \int \left(\int_0^t \phi(w_s, \theta) ds \right) \mathbb{B}x, x', t(dw). \end{aligned}$$

Finally, plugging this estimator into (B.3) yields an unbiased estimator of $\log \tilde{q}_t$.

References

- Aït-Sahalia, Y. (2008). Closed-form likelihood expansions for multivariate diffusions. *The Annals of Statistics*, 36(2):906–937.
- Beskos, A., Papaspiliopoulos, O., and Roberts, G. (2008). A factorisation of diffusion measure and finite sample path constructions. *Methodology and Computing in Applied Probability*, 10(1):85–104.
- Beskos, A., Papaspiliopoulos, O., Roberts, G., and Fearnhead, P. (2006). Exact and computationally efficient likelihood-based estimation for discretely observed diffusion processes. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 68(3):333–382. With discussions and a reply by the authors.
- Cappé, O., Moulines, E., and Rydén, T. (2005). *Inference in Hidden Markov Models*.
- Dempster, A., Laird, N., and Rubin, D. (1977). Maximum likelihood from incomplete data via the EM algorithm. *J. Roy. Statist. Soc. Ser. B*, 39:1–38.
- Douc, R., Fort, G., Moulines, E., and Priouret, P. (2009a). Forgetting the initial distribution for hidden markov models. *Stoch. Process. Appl.*, 119(4):1235–1256.
- Douc, R., Garivier, A., Moulines, E., and Olsson, J. (2009b). Sequential Monte Carlo smoothing for general state space hidden Markov models. Technical Report 2009-8, Lund University.
- Douc, R. and Moulines, E. (2008). Limit theorems for weighted samples with applications to sequential monte carlo methods. *Annals of Statistics*, 10.

- Doucet, A., Godsill, S., and Andrieu, C. (2000). On Sequential Monte Carlo Sampling Methods for Bayesian Filtering. *Statistics and Computing*, 10:197–208.
- Fearnhead, P., Papaspiliopoulos, O., and Roberts, G. (2008). Particle filters for partially observed diffusions. *Journal Of The Royal Statistical Society Series B*, 70(4):755–777.
- Fort, G. and Moulines, E. (2003). Convergence of the Monte Carlo expectation maximization for curved exponential families. *Ann. Stat.*, 31(4):1220–1259.
- Godsill, S. J., Doucet, A., and West, M. (2004). Monte carlo smoothing for nonlinear time series. *Journal of the American Statistical Association*, 99:156–168.
- Gordon, N., Salmond, D., and Smith, A. (1993). Novel approach to nonlinear/non-gaussian bayesian state estimation. *IEE Proc. F, Radar signal Process.*, 140:107–113.
- Handschin, J. and Mayne, D. (1969). Monte carlo techniques to estimate the conditional expectation in multi-stage non-linear filtering. *Int. J. Control*, 9:547–559.
- Ionides, E. L., Bhadra, A., and King, A. A. (2009). Iterated filtering. *arXiv:0902.0347*.
- Kitigawa, G. (1998). A self-organizing state-space-model. *Journal of the American Statistical Association*, 93(443):1203–1215.
- Kloeden, P. E. and Platen, E. (1992). *Numerical Solution of Stochastic Differential Equations*. Springer-Verlag, Berlin.
- Liu, J. and Chen, R. (1995). Blind deconvolution via sequential imputations. *J. Am. Statist. Assoc.*, 90(420):567–576.
- Olsson, J., Cappé, O., Douc, R., and Moulines, E. (2008). Sequential monte carlo smoothing with application to parameter estimation in nonlinear state space models. *Bernoulli*, 14(1):155–179.
- Olsson, J. and Rydén, T. (2008). Asymptotic properties of the bootstrap particle filter maximum likelihood estimator for state space models. *Stoch. Process. Appl.*, 118:649–680.
- Pedersen, A. R. (1995). Consistency and Asymptotic Normality of an Approximative Maximum Likelihood Estimator for Discretely Observed Diffusion Processes. *Bernoulli*, 1(3):257–279.
- Pitt, M. and Shephard, N. (1999). Filtering via simulation: Auxiliary particle filters. *J. Am. Statist. Assoc.*, 87:493–499.
- Wu, C. (1983). On the convergence properties of the EM algorithm. *Ann. Statist.*, 11:95–103.